

# articles

## Nucleotide sequence of bacteriophage $\Phi$ X174 DNA

F. Sanger, G. M. Air\*, B. G. Barrell, N. L. Brown†, A. R. Coulson, J. C. Fiddes, C. A. Hutchison III‡, P. M. Slocombe§ & M. Smith¶

MRC Laboratory of Molecular Biology, Hills Road, Cambridge CB2 2QH, UK

*A DNA sequence for the genome of bacteriophage  $\Phi$ X174 of approximately 5,375 nucleotides has been determined using the rapid and simple 'plus and minus' method. The sequence identifies many of the features responsible for the production of the proteins of the nine known genes of the organism, including initiation and termination sites for the proteins and RNAs. Two pairs of genes are coded by the same region of DNA using different reading frames.*

THE genome of bacteriophage  $\Phi$ X174 is a single-stranded, circular DNA of approximately 5,400 nucleotides coding for nine known proteins. The order of these genes, as determined by genetic techniques<sup>2-4</sup>, is *A-B-C-D-E-J-F-G-H*. Genes *F*, *G* and *H* code for structural proteins of the virus capsid, and gene *J* (as defined by sequence work) codes for a small basic protein that is also part of the virion. Gene *A* is required for double-stranded DNA replication and single-strand synthesis. Genes *B*, *C* and *D* are involved in the production of viral single-stranded DNA: however, the exact function of these gene products is not clear as they may either be involved directly in DNA synthesis or be required for DNA packaging, which is coupled with single-strand production. Gene *E* is responsible for lysis of the host.

The first nucleotide sequences established in  $\Phi$ X were pyrimidine tracts<sup>5-7</sup> obtained by the Burton and Petersen<sup>8</sup> depurination procedure. The longer tracts could be obtained pure and sequences of up to 10 nucleotides were obtained. More recently Chadwell<sup>9</sup> has improved the hydrazinolysis method to obtain the longer purine tracts. These results are included in the sequence given in Fig. 1.

More extensive  $\Phi$ X sequences were obtained using partial degradation techniques, particularly with endonuclease IV (refs 10 and 11). Ziff *et al.*<sup>12,13</sup> used this enzyme in conditions of partial hydrolysis to obtain fragments 50-200 nucleotides long which were purified as <sup>32</sup>P-labelled material by electrophoresis on polyacrylamide gels. The fragments came from the same region of the genome and the sequence of a 48-nucleotide long fragment (band 6, positions 1,047-1,094) was determined using mainly further degradation with endonuclease IV and partial exonuclease digestions.

Another 50-nucleotide long fragment was obtained by Robertson *et al.*<sup>14</sup> as a ribosome binding site. The viral (or plus)

strand DNA of  $\Phi$ X has the same sequence as the mRNA and, in certain conditions, will bind ribosomes so that a protected fragment can be isolated and sequenced. Only one major site was found. By comparison with the amino acid sequence data it was found that this ribosome binding site sequence coded for the initiation of the gene *G* protein<sup>15</sup> (positions 2,362-2,413).

At this stage sequencing techniques using primed synthesis with DNA polymerase were being developed<sup>16</sup> and Schott<sup>17</sup> synthesised a decanucleotide with a sequence complementary to part of the ribosome binding site. This was used to prime into the intercistronic region between the *F* and *G* genes, using DNA polymerase and <sup>32</sup>P-labelled triphosphates<sup>18</sup>. The ribo-substitution technique<sup>16</sup> facilitated the sequence determination of the labelled DNA produced. This decanucleotide-primed system was also used to develop the plus and minus method<sup>1</sup>. Suitable synthetic primers are, however, difficult to prepare and as DNA fragments generated by restriction enzymes are more readily available these have been used for most of the work reported here.

Another approach to DNA sequencing is to make an RNA copy using RNA polymerase with  $\alpha$ -<sup>32</sup>P-labelled ribotriphosphates and then to determine the RNA sequence by more established methods. Blackburn<sup>19,20</sup> used this approach on intact single-stranded  $\Phi$ X and on fragments obtained by digestion with endonuclease IV or with restriction enzymes. Sedat *et al.*<sup>21</sup> were extending their studies on the larger endonuclease IV fragments and their results, taken in conjunction with the transcription of the DNA fragments<sup>20</sup>, amino acid sequence of the *F* protein<sup>22</sup>, and the plus and minus method results, made it possible to deduce a sequence of 281 nucleotides (positions 1,016-1,296, Fig. 1) within the *F* gene<sup>23</sup>. Transcription of *Hind*II fragment 10, amino acid sequence data in the *G* gene, and the plus and minus method using *Hind*II fragments 2 and 10 as primers, gave a sequence of 195 nucleotides (positions 2,387-2,582, Fig. 1) at the N terminus of gene *G* (ref. 24).

### The 'plus and minus' method

Further work on the  $\Phi$ X sequence has been done using chiefly the plus and minus method primed with restriction fragments. Figure 2 shows the various restriction enzymes used and the fragment maps for each (refs 25-30 and C.A.H., submitted for publication, and N.L.B., C.A.H. and M.S., submitted for publication).

Figure 1 shows the combined results of the sequence work to date. The sequence is numbered from the single cleavage site of the restriction enzyme *Pst*I. As with other methods of sequencing nucleic acids, the plus and minus technique used by itself cannot be regarded as a completely reliable system and occasional errors may occur. Such errors and uncertainties can only

Present addresses: \*John Curtin School of Medical Research, Microbiology Department, Canberra City ACT 2601, Australia, †Department of Biochemistry, University of Bristol, Bristol BS8 1TD, UK, ‡Department of Bacteriology and Immunology, University of North Carolina, Chapel Hill, North Carolina 27514, §Max-Planck-Institut für Molekulare Genetik, 1 Berlin 33, FRG, ¶Department of Biochemistry, University of British Columbia, Vancouver BC, Canada V6T 1W5.

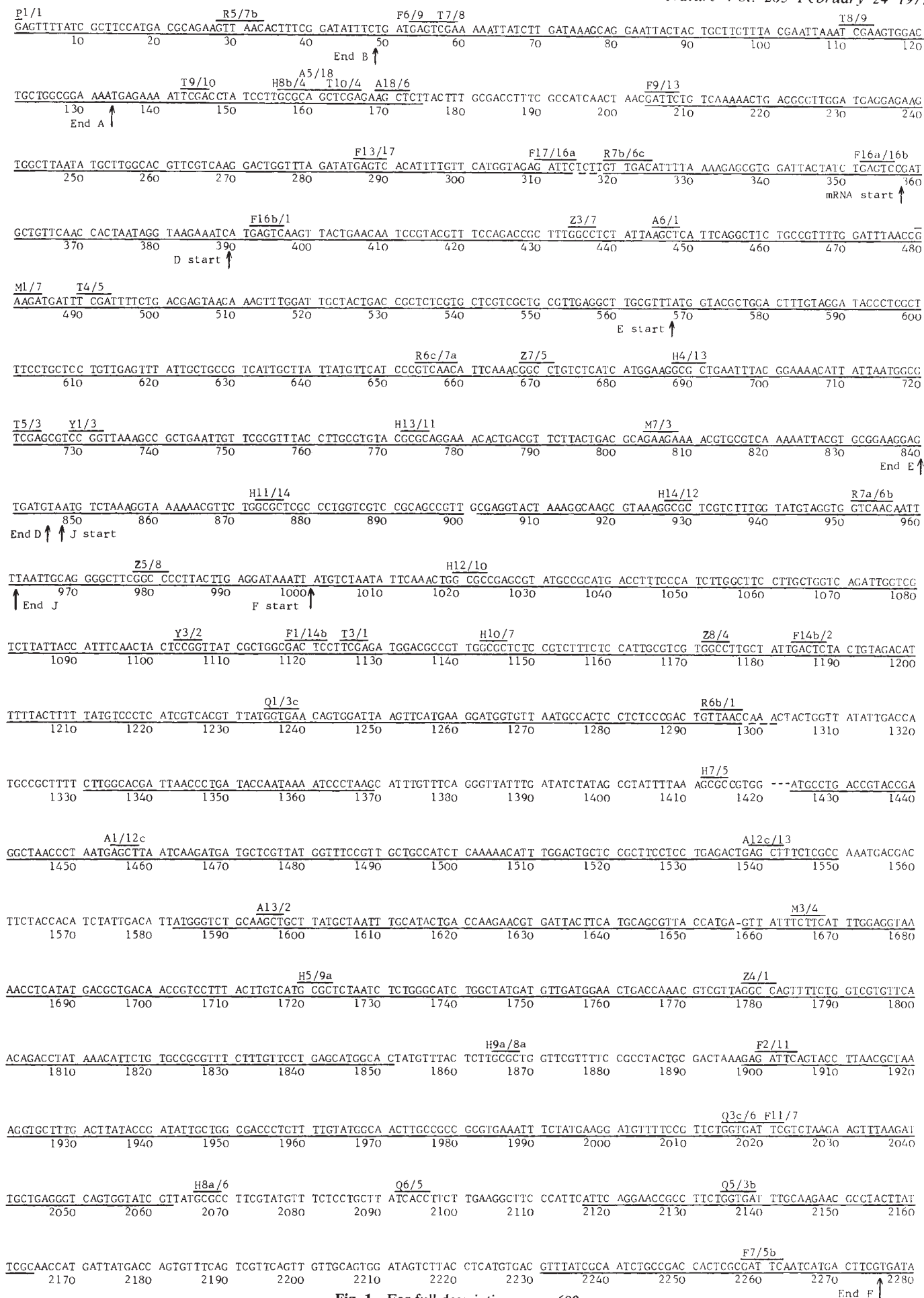


Fig. 1 For full description see p. 690.

AAAGATTGAG TGTCAGGTTA TAACCGAAGC GGTAAAAATT TTAATTTTIG CCGCTGAGGG GTGACCAAG CGAAGCCGGG TAGGTTTCT GCTTAGGAGT TTAATCATGT TTCAGACTTT  
 2290 2300 2310 2320 2330 2340 2350 2360 2370 2380 2390 2400  
 G start ↑

TAITTCCTGC CACAATCAA ACTTTTITTC TGATAAGCTG GTTCCTCACT CTGTACTCC AGCTTCTTCG GCACCTGTTT TACAGACACC TAAAGCTACA TCGTCAACGT TATATTTTGA  
 2410 2420 2430 2440 2450 2460 2470 2480 2490 2500 2510 2520  
 A2/16 A16/15a M4/10 A15a/3 R9/10

TAGTTTGAGG GTTAATGCTG GTAATGGTGG TITTCCTCAT TGCATTGAGA TGGATACATC TGTCACCGCC GCTAATCAGC TGTTCCTCAGT TGGTGTGAT ATTGCTTTTG ATGCCGACCC  
 2530 2540 2550 2560 2570 2580 2590 2600 2610 2620 2630 2640  
 M10/9 R10/2

TAAATTTTTT CCGCTGTTGG TTCGCTTTGA GTCTTCTTCG GTTCGGACTA CCCTCCCGAC TGCTTATGAT GTTTATCCIT TGGATGGTGG CCATGATGGT GTTATTATA CCGTCAAGGA  
 2650 2660 2670 2680 2690 2700 2710 2720 2730 2740 2750 2760  
 F5b/8 M9/2

CIGTIGACT ATTGAGGICC TCCCGCTAC GCCTGGCAAT AACGCTCAGC TGGTTCAT GGTTGGTCT AACCTTACCC CTACTAAATC CCGCGGATTC GTTCGGCTGA ATCAGGTTAT  
 2770 2780 2790 2800 2810 2820 2830 2840 2850 2860 2870 2880  
 Y2/5 F8/4

TAAAGAGATT ATTGTCTCC AGCCACTTAA GTGAGGTGAT TTATGTTTGG TGCTATTGCT GGCGGTAITG CTTCTGCTCT TGCTGGTGGC GCATGCTA AATTGTTTGG AGGCCGTCAA  
 2890 2900 2910 2920 2930 2940 2950 2960 2970 2980 2990 3000  
 End G ↑ H start ↑ Q3b/4 H3/2

AAAGCCGCT CCGGTGGCAT TCAAGGTGAT GTGCTTGCTA CCGATAACAA TACTGTAGGC ATGGGTGATG CTGGTATTA ATCTGCCATT CAAGGCTCTA ATCTCCCTAA CCCTGATGAG  
 3010 3020 3030 3040 3050 3060 3070 3080 3090 3100 3110 3120  
 Y5/4 Q4/7 Q7/2

GGCCGCCCTA GTTTCTTTC GTCTGCTATT GCTAAAGCTG GTAAGGACTT TCTTGAAGGT ACGTTGCAGG CTGGCACTTC TGCCCTTTC GATAAGTTCG TTGATTTGGT TGGACTTGGT  
 3130 3140 3150 3160 3170 3180 3190 3200 3210 3220 3230 3240  
 Z1/2 A3/9

GGCAAGTCTG CCGCTGATAA AGGAAAGGAT ACTCGTGATT ATCTGCTGCG TGCATTTCCT GAGCTTAAAG CTTGGGAGCG TGCTGGTGGT GATGCTTCCT CTGCTGGTAT GGTTCAGGCC  
 3250 3260 3270 3280 3290 3300 3310 3320 3330 3340 3350 3360  
 A9/12 R2/6a

GGATTTGGAG ATCAAAAAGA GCTTACTAAA ATGCAACICG ACAATCAGAA AGAGATTGCC GAGATGCAAA ATGAGACTCA AAAAGAGATT GCTGGCATTC AGTCGGGCAC TTCACGGCCAG  
 3370 3380 3390 3400 3410 3420 3430 3440 3450 3460 3470 3480  
 Y4/1 F4/14a A12d/7c F14a/12

AATACGAAAG ACCAGGTATA TGCACAAAAT GAGATGCTTG CTTATTC-AC AGAAGGAGTTC TACTGCTGCG TGCGTCTAT TATGGAAAAC ACCAATCTTT CCAAGCAACA GCAGGTTCCT  
 3490 3500 3510 3520 3530 3540 3550 3560 3570 3580 3590 3600  
 F12/10

GAGATTATGC GCCAAATGCT TACTCAAGCT CAAACGGCTG GTCAGTATTT TACCAATGAC CAAATCAAG AAATGACTCG CAAGGTTAGT GCTGAGTTG ACTTACTTCA TCAGCAAACC  
 3610 3620 3630 3640 3650 3660 3670 3680 3690 3700 3710 3720  
 H2/9b A7c/8 F10/15 R6a/4

CAGAATCAGC GGTATGGCTC TTCATATATT GGCGCTACTG CAAAGGATAT TTCATAITGC GTCACGTATG CTGCTCTGCG TGTGGTTGAT ATTTTTCATG GTATTGATAA AGCTGTTCCT  
 3730 3740 3750 3760 3770 3780 3790 3800 3810 3820 3830 3840  
 F15/5c M2/5 H9b/1 A8/14

GATACTTGGG ACAATTTCTG GAAAGACGGT AAAGCTGATG GTATTGGCTC TAATTGTCTT AGGAAATAAC CGTCAGGATT GACACCTTCC CAATGTATG TTTTCATGCC TCCAATCTT  
 3850 3860 3870 3880 3890 3900 3910 3920 3930 3940 3950 3960  
 End H ↑ mRNA start ↑ A14/7b

GGAGGCTTTT TTATGGTTCG TCTTATTAC CCTTCTGAAT GTCACGCTGA TTATTTTGGC TTTGGAGCTA TCGAGGCTCT TAAACCTGCT ATTGAGGCTT GTGGCATTTC TACTCTTCTT  
 3970 3980 3990 4000 4010 4020 4030 4040 4050 4060 4070 4080  
 A start ↑ T1/6

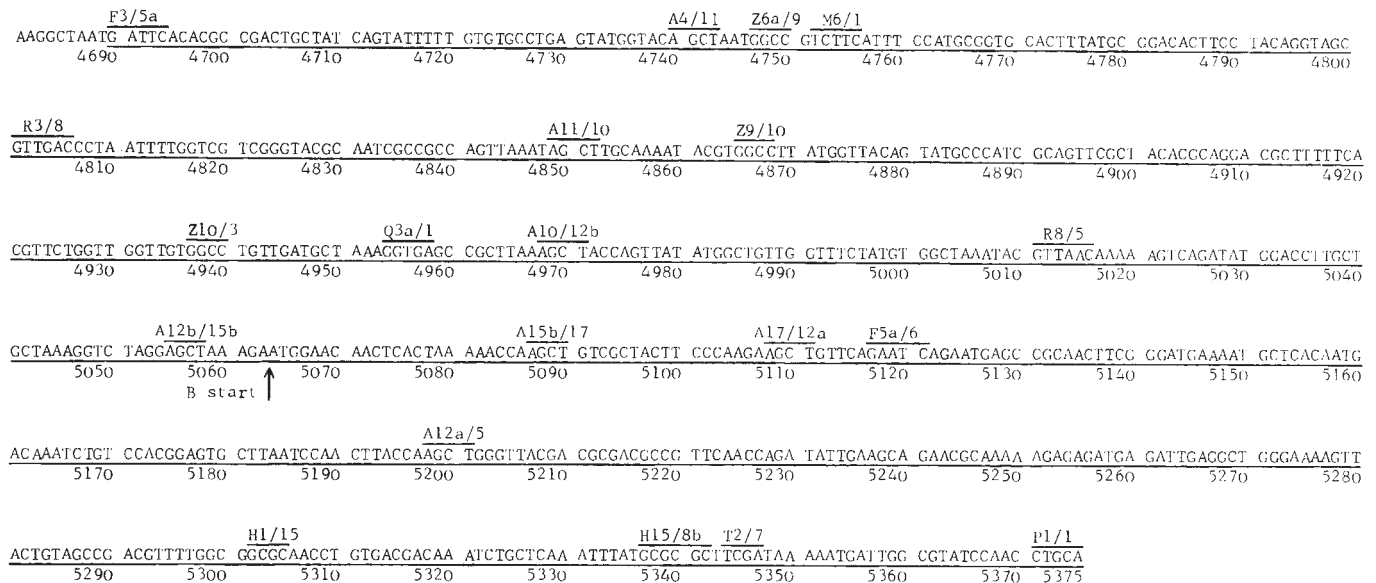
CAATCCCAA TGCTTGGCTT CCATAAGCAG ATGGATAACC GCATCAAGCT CTTGGAAAGAG ATCTCTGCTT TTCGATGCA GGCGTITGAG TTCGATAATG GTCATATGTA TGTTGACCCG  
 4090 4100 4110 4120 4130 4140 4150 4160 4170 4180 4190 4200  
 A7b/7a M5/8 F5c/3 T6/2 Q2/3a R4/3 Z2/6b

CAATAAGGCTG CTTCTGACGT TCGTGATGAG TTTGTATCTG TACTGAGAA GTTAATGGAT GAAATGGCAC AATGCTACAA TGTCCTCCCC CAACTGATA TTAATAACAC TATAGACACC  
 4210 4220 4230 4240 4250 4260 4270 4280 4290 4300 4310 4320

CGCCCCAAG GGGACGAAAA ATGGTTTTTA GAGAACCAGA AGACCGTTAC GCAGTTTTTG AAGCTGGCTG CTGAACGCC TCTTAAAGAT ATTCCGGAIG AGTATAAATA CCGCAAAAAA  
 4330 4340 4350 4360 4370 4380 4390 4400 4410 4420 4430 4440  
 M8/6 A7a/4

AAAGGTATTA AGGATGAGTG TTCAGAAATG CTGGAGGCC CCACTAAGAT ATCGCGTAGA GGCTTTGCTA TTCAGCGTTT GATGAATGCA ATCGGACAGG CTCATCTTGA TGGTTGGTTT  
 4450 4460 4470 4480 4490 4500 4510 4520 4530 4540 4550 4560  
 Z6b/6a

ATCGTTTTTG ACACTCTCAC GTGGCTGAC GACCGATTAG AGCGGTTTTA TGATAATCCC AATGCTTGGC GTGACTATIT TCGGATAAT GGTCTATGG TTCTTGTGTC CGACGGTGGC  
 4570 4580 4590 4600 4610 4620 4630 4640 4650 4660 4670 4680



**Fig. 1** A provisional nucleotide sequence for the DNA of bacteriophage  $\Phi$ X174 *am3* *cs70*. Solid underlining indicates sequences that are fully confirmed; sequences with no underlining probably do not contain more than one mistake per 50 residues. Broken underlining indicates more uncertain sequences. Restriction enzyme recognition sites are indicated (for key to single letter enzyme code see legend to Fig. 2), as are mRNA starts and protein initiation and termination sites. Nucleotides 4,127 to 4,201 have been independently sequenced by van Mansfield *et al.*<sup>58</sup>. The *am3* codon is at position 587.

be eliminated by more laborious experiments and, although much of the sequence has been so confirmed, it would probably be a long time before the complete sequence could be established. We are not certain that there is any scientific justification for establishing every detail and, as it is felt that the results may be useful to other workers, it has been decided to publish the sequence in its present form.

As template we have used both the viral (plus) and complementary (minus) strands of  $\Phi$ X. Usually it is possible to determine a sequence with a single primer starting at about 15–100 nucleotides from the appropriate restriction enzyme site. In a particularly good experiment the sequence can be read out to 150–200 nucleotides but the results may become less reliable. Most sequences have been derived by priming on both strands; this allows more confidence than when only one strand could be used.

A useful method for confirming runs of the same nucleotide is depurination of <sup>32</sup>P-labelled small restriction enzyme fragments or of products of the DNA polymerase priming experiments (ref. 31 and N.L.B. and M.S., in preparation). The most satisfactory way of confirming the DNA sequences is through amino acid sequence data. As the methods used are entirely unrelated, the results of the two approaches complement each other very well and therefore complete sequences can usually be deduced from incomplete data obtained by each method. The complete sequence of genes *G* (ref. 32), *D* (ref. 33), *J* (ref. 33 and Freymeyer, unpublished) and most of *F* have been obtained in this way.

Many of the sequences in Fig. 1 have been amply confirmed and are regarded as established: these are indicated in the figure by underlining. Some sequences are considered to be reasonably accurate and probably contain no more than one mistake in every 50 nucleotides. Sequences that are particularly uncertain—either because of lack of data or conflicting results—are also indicated in Fig. 1.

In considering the sequence of  $\Phi$ X174 as a functional unit it is convenient to begin in the region between the *H* and *A* genes and to continue around the DNA in the direction of transcription and translation.

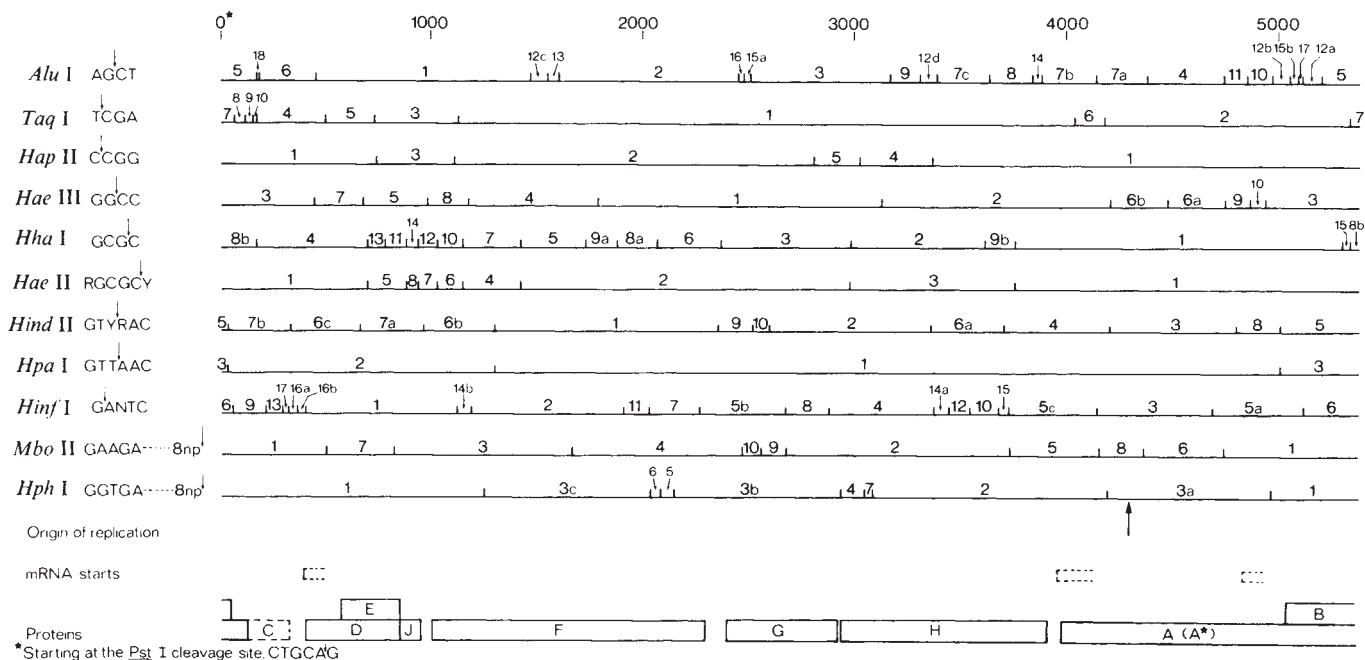
## *A* promoter and terminator

Sinsheimer *et al.*<sup>34,35</sup> and Axelrod<sup>36</sup> have determined the sequences of the 5' end of three  $\Phi$ X *in vitro* mRNA species and have located them on the restriction map. These sequences have been identified on the DNA sequence and one of them (AAATCTTGG) is found only at position 3,954 at which an *in vivo* unstable mRNA start has been located<sup>37</sup>. The sequence to the left of this has some characteristics of typical *E. coli* promoters<sup>38</sup> in that five out of the 'ideal' TATPuATPu residues are present. Nearby, to the right of this mRNA initiation, however, is the sequence TTTTITA which is similar to sequences found at the 3' ends of a number of mRNAs (see ref. 39) and seems a likely signal for mRNA termination. The presence of a rho-independent termination site in this approximate position has been suggested<sup>36,37</sup>, but the relative positions of the initiating and putative termination signals is rather surprising since the terminator for one mRNA would be expected to precede the initiator for the next. One possibility is that the T<sub>6</sub>A might be acting as an 'attenuator' involved in the control of mRNA production in a similar manner to that suggested for the tryptophan operon by Bertrand *et al.*<sup>40</sup>. If indeed it were acting as a transcription terminator one would expect a small RNA of 20 nucleotides to be produced, but no such product has yet been detected. Recent work, however, (Rosenberg, unpublished and ref. 41) indicates that termination may require the presence of a base-paired loop structure before the termination site. From the DNA sequence such a loop is probably present before the T<sub>6</sub>A sequence, but in mRNA starting from the initiation site at position 3,954 this loop is not formed (Fig. 3). Therefore mRNA that had started at an earlier promoter and extended through the *H* gene would be expected to terminate here, whereas mRNA newly initiated at position 3,954 would not. This could be a way in which the phage has economised on the use of DNA—by having the ends of the two mRNAs overlapping.

## The *A* protein

Where the amino acid sequence is available there is no problem in relating the DNA sequence to its coding properties, but it is





**Fig. 2** Fragment maps of restriction enzymes used in the sequence analysis of  $\Phi X174$  *am3* RFI DNA. Fragment maps of  $\Phi X174$  have been prepared for *Hind*II (R), *Hae*III (Z) and *Hpa*I+II by Lee and Sinshheimer<sup>25</sup>, *Hin*HI and *Hap*II (Y) by Hayashi and Hayashi<sup>26</sup>, and for *Alu*I (A) by Vereijken *et al.*<sup>27</sup> and for *Pst*I (P) by Brown and Smith<sup>30</sup>. B.G.B., G.M.A., C.A.H. and D. Jaffe prepared the *Hin*I (F) map, C.A.H. the *Hph*I (Q) map, and Jeppesen *et al.*<sup>28</sup> the *Hha*I (H), *Alu*I, *Hae*II and *Hap*II maps by using a rapid method depending on priming with DNA polymerase. A rapid two-dimensional hybridisation technique has been developed by C.A.H. (submitted for publication) and recently used for mapping *Mbo*II (M) (N.L.B., C.A.H., and M.S., submitted for publication) and *Taq*I (T)<sup>29</sup>. *Hha*I and *Hin*I maps have also been prepared by Baas *et al.*<sup>52</sup>.

more difficult to do so in the absence of such data, as is the case for the A protein. One way of identifying the reading phase of the DNA is from the distribution of nonsense codons. Over a sufficiently long sequence that is known to be coding for a single protein there is usually one phase that contains no nonsense codons, and this is identified as the reading phase. This requires completely accurate determination of the DNA sequences however: omission of a single nucleotide may give completely erroneous results. Another approach is possible in the case of  $\Phi X$ . The results with the *F* and *G* genes<sup>23, 24, 32</sup> showed an unexpectedly high frequency of codons ending in T. Therefore in a coding region there is a tendency for every third nucleotide to be a T and it is then possible to define the reading phase. Figure 4 illustrates how this characteristic was used to help determine the reading phase for the A protein and to identify its initiation codon at position 3,973. In a similar way the distribution of Ts may be used to identify errors in the DNA sequence, provided that such errors occur only infrequently.

A different approach to identifying the initiation site and reading phase in a coding sequence is by looking for a characteristic 'initiation sequence'. Shine and Dalgarno have shown that a common feature of ribosome binding sites is a number of nucleotides (at least three) preceding the ATG that are capable of forming base pairs with a sequence at the 3' end of 16S rRNA<sup>42, 43</sup>. All of the known initiation sites in  $\Phi X174$  that have been identified by direct amino acid sequencing (for the F, G, H, J and D proteins) satisfy this criterion (see Table 2) and the fact that the sequence preceding the ATG in position 3,973 also has this characteristic supports its identification as the initiation site for the A protein.

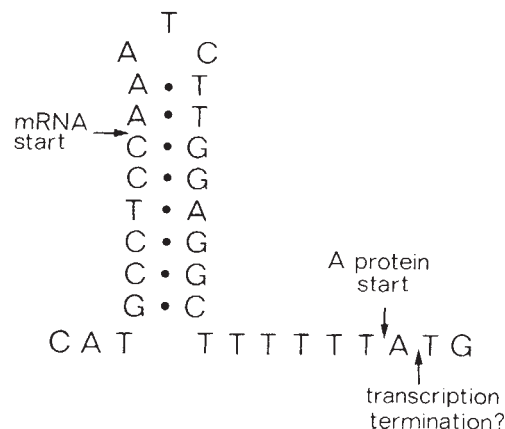
If, as has been suggested<sup>37</sup>, some mRNA from the previous promoter does extend beyond the hairpin structure, initiation of A protein synthesis may be controlled by the inclusion of the region complementary to the 16S rRNA in the hairpin loop. This could explain the presence of two types of mRNA covering the *A* cistron, as suggested by Hayashi *et al.*<sup>37</sup>—one unstable and active and the other stable but inactive. The former would be initiated at the A promoter, the latter at an earlier promoter and result from 'read-through' at the terminator. The postu-

lated reading frame for the A protein was confirmed by sequencing amber mutants mapping in the N-terminal region of gene *A*. *am86* proved to be a C→T change at position 4,108 and *am33* a C→T change at position 4,372. These both result in formation of an amber codon (TAG) in the same reading frame as the proposed initiating ATG and the sequence continues to the termination codon at position 133. The A protein, which is the largest coded by  $\Phi X174$ , is thus 512 amino acids long with a molecular weight of 56,000, in good agreement with SDS gel estimations (see refs 4 and 44). The A\* protein, with a molecular weight of about 35,000, is believed to result from an internal translational start in the *A* gene, in the same reading phase<sup>45</sup>. From consideration of possible ribosome binding sequences<sup>42, 43</sup> the ATG in position 4,657 seems to be the most likely initiation site for the A\* protein.

**The origin of replication**

The origin of  $\Phi X$  viral strand DNA synthesis has been located in gene *A*, in restriction fragment Z6b (ref. 46). This origin, while

**Fig. 3** Potential secondary structure at the *A* mRNA start.



	Ts in codon position		
	(1)	(2)	(3)
3,910 CCG . TCA . GGA . TTG . ACA . CCC . TCC . CAA . TTG . TAT .	5	2	1
3,940 GTT . TTC . ATG . CCT . CCA . AAT . CTT . GGA . GGC . TTT .	3	4	3*
3,970 TTT . <u>ATG</u> . GTT . CCT . TCT . TAT . TAC . CCT . TCT . GAA .	2	5	5
4,000 TGT . CAC . GCT . GAT . TAT . TTT . GAC . TTT . GAG . CGT .	3	5	7*
	5	0	7*
	4	2	7

\*These figures refer to the last five codons of the previous line and the first five of the next line.

coding for part of the A protein, probably corresponds to the position of the plus strand nick made by the same protein<sup>44</sup>. Gaps in this region that are found in replicating double-stranded (RF) DNAs are probably related to the position of the nick. Eisenberg *et al.*<sup>47</sup> have investigated such gaps by depurination analysis and identified, in particular, the product C<sub>6</sub>T. The sequence CTC<sub>5</sub> is found in position 4,285 (Fig. 1) and the location of the origin in this region agrees precisely with the results of Baas *et al.*<sup>46</sup>. It is not possible at present to identify the actual position of the origin nick. The region shows no apparent secondary structure or symmetrical sequences, although there is an AT-rich region (4,298–4,307) between two GC-rich regions which might be of significance. Such a region is found near the origin of replication of SV40 DNA (ref. 48).

### B promoter

The second of the mRNA 5' sequences (AUCGC)<sup>34</sup> has been mapped in restriction fragment R8 (Fig. 2), which starts about 300 nucleotides on from the proposed A\* initiation. The sequence ATCGC is found at positions 4,832 and 4,888 in Fig. 1. The only way we can choose between them at the moment is that the second is preceded by the sequence TACAGTA (position 4,877), which is more akin to sequences found in known promoters<sup>38</sup> than are sequences preceding the other possible mRNA start. Irrespective of which of these sequences is used, the mRNA has a long 'leader' sequence (232

or 176 nucleotides) before the next proposed initiation codon (gene B).

### The B protein

From a study of the ribonuclease T<sub>1</sub> digestion products of the ribosome binding sites of ΦX mRNAs<sup>49</sup>, it was possible to identify an initiating ATG in position 5,064. From the genetic map<sup>2,3</sup>, this would be expected to be gene B but, as discussed above, the A protein coding sequence extends right through this region, past the *Pst* site at residue 1 in Fig. 1, and terminates at residue 133. The initiating codon contained in the ribosome-protected sequence is, however, out of phase with the A protein reading frame. The proposed B protein coding sequence is one nucleotide to the left of the A protein phase, and continues until a termination codon occurs at position 49. Therefore the B protein coding sequence is totally contained within the A gene. These reading frames have been confirmed by sequencing mutants in genes A and B (*am16*, N.L.B. and M.S., in preparation; *am18*, *am35*, *ts116* (ref. 50)). Since the B protein has not been purified no protein sequence data is available. The complete amino acid sequence can be predicted from the DNA sequence however. The protein is 120 amino acids long with a molecular weight of 13,845 (including the N-terminal Met). The molecular weight estimates of the B protein obtained by SDS-gel electrophoresis are mostly greater than this (see review, ref. 4), but the electrophoretic mobility varied with gel concentration and cross linker. Such anomalous behaviour suggests that there may be, for instance, carbohydrate attached to the B protein.

### The C protein

The next known gene product, protein C, maps between genes B and D. Examination of the DNA sequence in this region indicates that the most probable initiating ATG overlaps the termination codon, TGA, of gene A in the sequence ATGA at position 134. A possible termination codon for gene C could then be at position 391, although the sequence and phasing is not yet confirmed through this region. There is another possible protein initiation codon (position 51, overlapping the B protein termination codon) which would result in a slightly shorter gene product terminating at nucleotide 219. For the C protein, however, we favour the 'A terminator' start, since only this reading frame contains a CAA sequence, which by a C → T alteration could give the ochre 6 mutant. Ochre 6 is a gene C mutant produced by the decay of <sup>3</sup>H-cytosine<sup>51</sup> and has been mapped in fragments A6 and F9 (ref. 52); that is, between nucleotides 170 and 205 (Fig. 1).

### Sequence following the D promoter

The mRNA 5' sequence which maps before the D gene (GAUGC)<sup>34</sup> is found at position 358 in Fig. 1. The sequence preceding the messenger start has only four of the TATPuATPu nucleotides<sup>38</sup>. Thirty-two nucleotides after the mRNA initiation is the ATG (position 390) that initiates D protein synthesis. The amino acid sequence of the D protein has been determined

Table 1 ΦX174 coding capacity

Gene	Protein molecular weight from SDS gels*	Number of nucleotides (Fig. 1)	Protein molecular weight from sequence information
A	55,000–67,000	1,536	56,000
(A*)	35,000		
B	19,000–25,000	(360)†	13,845‡
C	7,000		
D	14,500	456	16,811‡
E	10,000–17,500	(273)†	9,940
J	5,000	114	4,097‡
F	48,000	1,275	46,400
G	19,000	525	19,053‡
H	37,000	984	35,800
Non-coding and C		485	
Total		5,375	

\*See ref. 4.

†Values in parenthesis are overlapping sequences and therefore not included in the addition to obtain the total length of DNA.

‡These values are calculated from the amino acid sequence (in the case of B deduced from the nucleotide sequence). The others are derived using the formula

$$\text{Protein molecular weight} = \frac{\text{No. of nucleotides}}{3 \times 0.00915}$$

**Table 2** Initiation sequences of  $\Phi$ X174 coded proteins

D	C-C-A-C-T- <u>A-A-T</u> - <u>A-G-G-T</u> -A-A-G-A-A-A-T-C-A- <u>T-G-A-G-T</u> -C-A-A-G-T-T-A-C-T	Ser Gln Val Thr
E	C-T-G-C-G-T-T- <u>G-A-G-G</u> -C-T-T-G-C-G-T-T-T-A- <u>T-G-G-T</u> -A-C-G-C-T-G-G-A-C-T	
J	C-G-T-G-C-G-G- <u>A-A-G-G-A-G</u> - <u>T-G-A-T</u> -G-T-A- <u>A-T-G-T</u> -C-T-T-A-A-A-G-G-T-A-A-A	Ser Lys Gly Lys
F	C-C-C-T-T-A-C-T-T-G- <u>A-G-G-A</u> -T-T-A-A-T-T-A- <u>T-G-T</u> -C-T-A-A-T-A-T-T-C-A-A	Ser Asn Ile Gln
G	T-T-C-T-G-C-T-T- <u>A-G-G-A-G</u> -T-T-T-A-A-T-C- <u>A-T-G-T</u> -T-T-C-A-G-A-C-T-T-T-T	Met Phe Gln Thr Phe
H	C-C-A-C-T- <u>T-A-A-G</u> - <u>T-G-A-G-G-T-G-A-T</u> -T-T-A- <u>T-G-T</u> -T-T-G-G-T-G-C-T-A-T-T	Met Phe Gly Ala Ile
A	C-A-A-A-T-C-T-T- <u>G-G-A-G-G</u> -C-T-T-T-T-T-T-A- <u>T-G-G-T</u> -T-C-G-T-T-C-T-T-A-T	
B	A-A-A-C-G-T-C-T- <u>A-G-G-A-G</u> -C-T-T-A-A-A-G-A- <u>A-T-G-G</u> -A-A-C-A-A-C-T-C-A-C-T	
16S RNA		
3' end	HO <sup>A-U-U-C-C-U-C-C-A-C-U-A-G</sup>	

Where the protein start has been independently confirmed by protein sequencing data the amino acid sequences are indicated. The other initiation regions were identified as described in the text. Sequences complementary to the 3' end of 16S rRNA (refs 42, 43) are boxed; broken lines indicate further complementarity if some nucleotides are looped out or not matched. Ribosome binding to mRNA has been demonstrated in these regions for genes *J*, *F*, *G* and *B* (ref. 49).

almost completely, and nucleotide and amino acid sequences can be correlated to the termination codon at position 846 (ref. 33). The D protein, which is involved in capsid assembly, is 151 amino acids in length, with a molecular weight of 16,811. The D termination codon overlaps the initiation codon for gene *J* in the sequence TAATG. A similar structure has also been found by Platt and Yanofsky<sup>53</sup> in the tryptophan operon. The DNA sequence following this initiation codon matches the amino acid

sequence of the small basic protein (37 amino acids) of the virion determined by D. Freymeyer, P. R. Shank, T. Vanaman, C.A.H. and M. H. Edgell (personal communication). Benbow *et al.*<sup>2,3</sup> suggested that the mutation *am6* was located in a gene *J*, coding for the small protein of the virion, and mapping immediately before gene *F*. Although marker rescue experiments indicate that *am6* is not in this region<sup>54</sup>, the DNA sequence shows that there is a gene coding for the virion protein and we

**Table 3** Promoter sequences in  $\Phi$ X174

A promoter	A-G-G-A-T-T-G-A-C-A-C-C-C-T-C-C-C-A-A-T-T-G-T-A-T-G-T- <u>T-T-T-C-A-T-G</u> -C-C-T-C-C- <u>A-A-A-T-C-T</u> ...	3954 ↑ 18 nucleotides to A protein start
D promoter	<u>R7b/R6c</u> G-T-T-G-A-C-A-T-T-T-T-A-A-A-A-G-A-G-C-G-T-G-G-A-T-T-A-C- <u>T-A-T-C-T-G-A</u> -G-T-C-C- <u>G-A-T-G-C-T</u>	358 ↑ 32 nucleotides to D protein start
B promoter?	<u>R3/R8</u> C-A-G-G-T-A-G-C-G-T-T-G-A-C-C-C-T-A-A-T-T-T-T-G-G-T-C-G- <u>T-C-G-G-G-T-A</u> -C-G-C-A- <u>A-T-C-G-C-C</u>	4832 ↑ 232 nucleotides to B protein start
B promoter?	A-G-C-T-T-G-C-A-A-A-A-T-A-C-G-T-G-G-C-C-T-T-A-T-G-G-T- <u>T-A-C-A-G-T-A</u> -T-G-C-C- <u>C-A-T-C-G-C-A</u>	4888 ↑ 176 nucleotides to B protein start

mRNA initiation sequences<sup>34-36</sup> are underlined. Boxed regions indicate sequences that may correspond to the TATPuATPu sequence found in other promoters<sup>55</sup>, taking into account the distance from the mRNA starts.

have defined this as gene *J* (ref. 33). Since the *J* initiation codon overlaps the *D* termination codon we had to look elsewhere for gene *E*, which genetic mapping<sup>2,3</sup> had placed between them. Amber mutants in gene *E* (*am3*, *am27*, *am34* and *amN11*) were located by the marker rescue technique and sequenced. All were found to be within the *D* coding sequence, with the mutant amber codons one nucleotide to the right of the *D* reading frame<sup>33</sup>. Thus the *E* coding sequence is completely contained within the *D* coding region but in a different reading frame. The proposed initiation and termination codons for the *E* protein are at nucleotides 568 and 840, respectively<sup>33</sup>, giving a protein 91 amino acids in length with a molecular weight of about 9,900 (including the N-terminal methionine).

### The F protein

Following the *J* gene is an intercistronic region of 39 nucleotides before initiation of the *F* protein. There is no known function of this apparently untranslated sequence, although the presence of a hairpin structure (positions 969–984) suggests that it could be the site of the *in vivo* messenger termination signal<sup>37</sup> mapped in this region. The *F* protein is initiated by the ATG at position 1,001. This is the capsid component of the virion, and almost all the amino acid sequence is known<sup>22,24</sup>. There are regions in this gene where the DNA sequence is not completely established, but the protein is about 424 amino acids in length, giving a molecular weight of  $\approx 46,300$ .

### The G protein region

The termination signal for the *F* protein (position 2,276) is followed by an unusually long untranslated sequence of 111 nucleotides until the *G* protein initiation codon<sup>31</sup>. This region contains a looped structure which was postulated to have some functional role, as yet unknown, in the single-stranded DNA or the mRNA.

Initiation of the *G* protein at position 2,387 is followed by a sequence of 425 nucleotides until termination at position 2,912, giving a spike protein of molecular weight 19,053. The nucleotide and amino acid sequences of this gene and product are known<sup>24,32</sup>.

### The H protein

The initiation codon for the *H* protein (position 2,923) was identified first on the basis of the distribution of T nucleotides between the three reading phases, and later confirmed by amino acid sequence analysis. Amino acid sequence data on the *H* protein is minimal but the five peptide sequences known do correspond to the amino acid sequence, deduced from the DNA sequence by using the high frequency of third position T to help in assigning a reading frame to any given region. The DNA sequence is not entirely confirmed but it is possible to write a

reasonably accurate amino acid sequence for the *H* protein. The protein terminates at nucleotide 3,907, in agreement with carboxypeptidase results, giving a spike protein of molecular weight  $\approx 35,600$  (326 amino acids). The amino acid sequence at the N terminus seems to be particularly rich in hydrophobic residues, which is consistent with its suggested function as the 'pilot' protein that reacts with the bacterial membrane<sup>55,56</sup>. After *H* protein termination there are 66 nucleotides before initiation of the *A* protein at position 3,973.

### Coding capacity of the $\Phi$ X174 genome

The most striking feature of the  $\Phi$ X DNA sequence is the way in which the various functions of the genome are compressed within the 5,375 nucleotides. Since the identification of  $\Phi$ X gene products<sup>2,4</sup> it has been clear that proteins of the accepted molecular weights could not be separately coded on the available length of DNA. However, with the presence of two pairs of overlapping genes (*B* within *A* (ref. 50), *E* within *D* (ref. 33)) the genome has more coding capacity than had been originally supposed on the assumption that each gene was physically separate. Table 1 summarises the molecular weights of the known  $\Phi$ X-coded proteins. There are other potential initiation sites for polypeptide synthesis (for example, in genes *A*, *F*, *G* and *H*) and further genetic work may clarify whether there are in fact other  $\Phi$ X genes as yet unidentified.

### Initiation of protein synthesis

Table 2 lists the protein initiation sequences for genes *A*, *B*, *D*, *E*, *J*, *F*, *G* and *H*. It can be noted that there are no extra precursor sequences in proteins *D*, *J*, *F*, *G* or *H* at either the N or C terminus. There seems to be no relationship between the degree of complementarity to the 16S rRNA and the amount of protein synthesised, and we see no other features in the sequence that could explain different efficiencies of translation except where genes overlap.

### Transcription of $\Phi$ X174

The sequences preceding known mRNA starts<sup>34–36</sup> are shown in Table 3. Other studies on promoter sequences<sup>38</sup> have suggested certain features that they may have in common. Although some of these features are present in the sequences preceding the  $\Phi$ X transcription initiations others are not, and at present it is difficult to suggest what signal on the DNA determines a promoter site or the efficiency with which it initiates RNA synthesis. It is interesting to note that a polymerase binding site found by Chen *et al.*<sup>37</sup>, but not associated with any *in vitro* or *in vivo* mRNA starts, mapped near the region where there is the sequence TATGATG characteristic of promoters<sup>38</sup> (positions 2,705–2,711).

Table 4 Codons used in  $\Phi$ X174

Phe	TTT	39	Ser	TCT	35	Tyr	TAT	36	Cys	TGT	12
	TTC	26		TCC	9		TAC	15		TGC	10
Leu	TTA	19		TCA	16	Ter	TAA	3	Ter	TGA	5
	TTG	26		TCG	14		TAG	0	Trp	TGG	16
Leu	CTT	36	Pro	CCT	34	His	CAT	16	Arg	CGT	40
	CTC	15		CCC	6		CAC	7		CGC	29
	CTA	3		CCA	6	Gln	CAA	27		CGA	4
	CTG	24		CCG	21		CAG	34		CGG	8
Ile	ATT	45	Thr	ACT	40	Asn	AAT	37	Ser	AGT	9
	ATC	12		ACC	18		AAC	25		AGC	5
	ATA	2		ACA	13	Lys	AAA	47	Arg	AGA	6
Met	ATG	42		ACG	19		AAG	31		AGG	1
Val	GTT	53	Ala	GCT	64	Asp	GAT	44	Gly	GGT	38
	GTC	14		GCC	17		GAC	35		GGC	28
	GTA	10		GCA	12	Glu	GAA	27		GGA	13
	GTG	11		GCG	12		GAG	34		GGG	3

The totals are derived from sequences in Fig. 1 which are fully confirmed, that is, 377 codons in gene *A*, 120 in gene *B*, 152 in gene *D*, 91 in gene *E*, 38 in gene *J*, 344 in gene *F*, 175 in gene *G* and 49 in gene *H*. Out of a total of 1,346 codons 42.9% terminate in T. The percentages in the different genes are: *A*, 37.1 (non-overlapping region 47.1; overlapping region 15.8); *B*, 34.2; *D*, 42.1; *E*, 14.3; *J*, 47.4; *F*, 52.0; *G*, 54.3; *H*, 49.0. The initiating ATG is included in all cases.



## The use of codons in $\Phi X174$

Table 4 shows the codons used in regions where the nucleotide sequence is fully confirmed. It is clear that the pattern established by early observations on non-random use of codons<sup>23,24</sup> is continued now that more information is available. In particular, the preference for T at the third position of the codon is marked throughout the genome, as shown in Table 4. In regions of overlapping genes, one of the pair tends to continue the 'third T' trend (*D* and *B*), thus excluding the other (*E* and *A*). This may give some indication of the order in which overlapping genes evolved<sup>33,50</sup>. Another interesting feature is the very low occurrence of codons starting AG, particularly in non-overlapping regions. The base composition of the sequence of  $\Phi X174$  DNA shown in Fig. 1 is: A, 23.9%; C, 21.5%; G, 23.3% and T, 31.2%. This is in good agreement with previously determined values (see ref. 4).

We thank D. McCallum and R. Staden for carrying out the computer data storage and analysis of the sequence.

*Note added in proof:* J. E. Sims and D. Dressler (personal communication) have independently determined the sequence in positions 263–375 and 4,801–4,940. Their results agree with those given in Fig. 1. They have also identified the 'B' mRNA start as being at position 4,888.

Received November 30; accepted December 24, 1976.

- 1 Sanger, F. & Coulson, A. R. *J. molec. Biol.* **94**, 441–448 (1975).
- 2 Benbow, R. M., Hutchison, C. A. III, Fabricant, J. D. & Sinsheimer, R. L. *J. Virol.* **7**, 549–558 (1971).
- 3 Benbow, R. M., Zuccarelli, A. J., Davis, G. C. & Sinsheimer, R. L. *J. Virol.* **13**, 898–907 (1974).
- 4 Denhardt, D. T. *CRC Crit. Rev. Microbiol.* **4**, 161–222 (1975).
- 5 Hall, J. B. & Sinsheimer, R. L. *J. molec. Biol.* **6**, 115–127 (1963).
- 6 Ling, V. *Proc. natn. Acad. Sci. U.S.A.* **69**, 742–746 (1972).
- 7 Harbers, B., Delancy, A. D., Harbers, K. & Spencer, J. H. *Biochemistry* **15**, 407–414 (1976).
- 8 Burton, K. & Petersen, G. B. *Biochem. J.* **75**, 17–27 (1960).
- 9 Chadwell, H. A. Thesis, University of Cambridge (1974).
- 10 Sadowski, P. D. & Baktya, I. *J. biol. Chem.* **247**, 405–412 (1972).
- 11 Ling, V. *FEBS Lett.* **19**, 50–54 (1971).
- 12 Ziff, E. B., Sedat, J. W. & Galibert, F. *Nature new Biol.* **241**, 34–37 (1973).
- 13 Galibert, F., Sedat, J. W. & Ziff, E. B. *J. molec. Biol.* **87**, 377–407 (1974).

- 14 Robertson, H. D., Barrell, B. G., Weith, H. L. & Donelson, J. E. *Nature new Biol.* **241**, 38–40 (1973).
- 15 Air, G. M. & Bridgen, J. *Nature new Biol.* **241**, 40–41 (1973).
- 16 Sanger, F., Donelson, J. E., Coulson, A. R., Kössel, H. & Fischer, D. *Proc. natn. Acad. Sci. U.S.A.* **70**, 1209–1213 (1973).
- 17 Schott, H. *Makromolek. Chem.* **175**, 1683–1693 (1974).
- 18 Donelson, J. E., Barrell, B. G., Weith, H. L., Kössel, H. & Schott, H. *Eur. J. Biochem.* **58**, 383–395 (1975).
- 19 Blackburn, E. H. *J. molec. Biol.* **93**, 367–374 (1975).
- 20 Blackburn, E. H. *J. molec. Biol.* **107**, 417–432 (1976).
- 21 Sedat, J. W., Ziff, E. B. & Galibert, F. *J. molec. Biol.* **107**, 391–416 (1976).
- 22 Air, G. M. *J. molec. Biol.* **107**, 433–444 (1976).
- 23 Air, G. M. *et al. J. molec. Biol.* **107**, 445–458 (1976).
- 24 Air, G. M., Blackburn, E. H., Sanger, F. & Coulson, A. R. *J. molec. Biol.* **96**, 703–719 (1975).
- 25 Lee, A. S. & Sinsheimer, R. L. *Proc. natn. Acad. Sci. U.S.A.* **71**, 2882–2886 (1974).
- 26 Hayashi, M. N. & Hayashi, M. *J. Virol.* **14**, 1142–1152 (1974).
- 27 Vereijken, J. M., van Mansfeld, A. D. M., Baas, P. D. & Jansz, H. S. *Virology* **68**, 221–233 (1975).
- 28 Jeppesen, P. G. N., Sanders, L. & Slocombe, P. M. *Nucl. Acids Res.* **3**, 1323–1339 (1976).
- 29 Sato, S., Hutchison, C. A. III & Harris, J. I. *Proc. natn. Acad. Sci. U.S.A.* (in the press).
- 30 Brown, N. L. & Smith, M. *FEBS Lett.* **65**, 284–287 (1976).
- 31 Fiddes, J. C. *J. molec. Biol.* **107**, 1–24 (1976).
- 32 Air, G. M., Sanger, F. & Coulson, A. R. *J. molec. Biol.* **108**, 519–533 (1976).
- 33 Barrell, B. G., Air, G. M. & Hutchison, C. A. III *Nature* **264**, 34–41 (1976).
- 34 Smith, L. H. & Sinsheimer, R. L. *J. molec. Biol.* **103**, 699–735 (1976).
- 35 Grohmann, K., Smith, L. H. & Sinsheimer, R. L. *Biochemistry* **14**, 1951–1955 (1975).
- 36 Axelrod, N. *J. molec. Biol.* **108**, 753–779 (1976).
- 37 Hayashi, M., Fujimura, F. K. & Hayashi, M. *Proc. natn. Acad. Sci. U.S.A.* **73**, 3519–3523 (1976).
- 38 Pribnow, D. *Proc. natn. Acad. Sci. U.S.A.* **72**, 784–788 (1975).
- 39 Rosenberg, M., de Crombrughe, B. & Musso, R. *Proc. natn. Acad. Sci. U.S.A.* **73**, 717–721 (1976).
- 40 Bertrand, K. *et al. Science* **189**, 22–26 (1975).
- 41 Sugimoto, K., Sugisaki, H., Okamoto, T. & Takanami, M. *J. molec. Biol.* (in the press).
- 42 Shine, J. & Dalgarno, L. *Proc. natn. Acad. Sci. U.S.A.* **71**, 1342–1346 (1974).
- 43 Steitz, J. A. & Jakes, K. *Proc. natn. Acad. Sci. U.S.A.* **72**, 4734–4738 (1975).
- 44 Henry, T. J. & Knippers, R. *Proc. natn. Acad. Sci. U.S.A.* **71**, 1549–1553 (1974).
- 45 Linney, E. & Hayashi, M. *Nature* **249**, 345–348 (1974).
- 46 Baas, P. D., Jansz, H. S. & Sinsheimer, R. L. *J. molec. Biol.* **102**, 633–656 (1976).
- 47 Eisenberg, S., Harbers, B., Hours, C. & Denhardt, D. T. *J. molec. Biol.* **99**, 107–123 (1975).
- 48 Subramanian, K. N., Dhar, R. & Weissman, S. M. *J. biol. Chem.* (in the press).
- 49 Ravech, J. V., Model, P. & Robertson, H. D. *Nature* **265**, 698–702 (1977).
- 50 Smith, M. *et al.* (submitted to Nature).
- 51 Funk, F. & Sinsheimer, R. L. *J. Virol.* **6**, 12–19 (1970).
- 52 Baas, P. D., van Heusden, G. P. H., Vereijken, J. M., Weisbeek, P. J. & Jansz, H. S. *Nucl. Acids Res.* **3**, 1947–1960 (1976).
- 53 Platt, T. & Yanofsky, C. *Proc. natn. Acad. Sci. U.S.A.* **72**, 2399–2403 (1975).
- 54 Weisbeek, P. J., Vereijken, J. M., Baas, P. D., Jansz, H. S. & Van Arkel, G. A. *Virology* **72**, 61–71 (1976).
- 55 Jazwinski, S. M., Lindberg, A. A. & Kornberg, A. *Virology* **66**, 283–293 (1975).
- 56 Kornberg, A. *DNA Synthesis* (W. H. Freeman, San Francisco, 1974).
- 57 Chen, C. Y., Hutchison, C. A. III & Edgell, M. H. *Nature new Biol.* **243**, 233–236 (1973).
- 58 van Mansfeld, A. D. M., Vereijken, J. M. & Jansz, H. S. *Nucl. Acids Res.* **3**, 2827–2843 (1976).

# DNA sequence of a region of the $\Phi X174$ genome coding for a ribosome binding site

Nigel L. Brown\* & Michael Smith†

MRC Laboratory of Molecular Biology, Hills Road, Cambridge CB2 2QH, UK

*The DNA region corresponding to a newly identified ribosome binding site in  $\Phi X174$  DNA is sequenced and mapped in relation to the physical map of  $\Phi X174$ . Assignments of the binding site to a specific gene are discussed, and the possibility of a second case of overlapping genes in  $\Phi X174$  is considered.*

RECENT studies on the DNA of bacteriophage  $\Phi X174$  have defined sequences corresponding to the ribosome binding sites of genes *D*, *J*, *F* and *G* (refs 1–3 and B. G. Barrell, personal communication). A ribosome binding site that does not correspond to any of the above sequences has now been defined by studies on a ribosome-protected fragment of a transcript *in vitro* of  $\Phi X174$  RF DNA (ref. 4). Inspection of preliminary DNA sequence data obtained from  $\Phi X174$  (F. Sanger, *et al.*, unpublished) for potential ribosome binding sites, and comparison of the sequence data with T<sub>1</sub>-oligonucleotides from the unassigned RNA fragment<sup>4</sup>, allowed identification of the

DNA sequence corresponding to this RNA fragment. We report here the complete DNA sequence of a region of  $\Phi X174$  *am3* RF DNA that corresponds to this new ribosome binding site. Because the DNA sequence contains accurately mapped cleavage sites for restriction endonucleases, the ribosome binding site can be located precisely on the physical map of the  $\Phi X$  genome<sup>5</sup>, and therefore located on the genetic map<sup>6–8</sup>.

## Determination of the DNA sequence

The ribosome binding site lies in that portion of *Hind*II fragment 5 that overlaps *Hinf*I fragment 5a (nucleotides 5,000–5,100 approx, Fig. 1). The DNA sequence of this region was determined using the 'plus-minus' technique of Sanger and Coulson<sup>9</sup> (Fig. 2). It was quantitatively confirmed by characterisation of the pyrimidine oligonucleotides obtained by primed synthesis *in vitro* of DNA labelled with  $\alpha$ -<sup>32</sup>P-dATP or  $\alpha$ -<sup>32</sup>P-dGTP using viral or complementary strand as template<sup>10,11</sup> (Fig. 3, Table 1).

The sequence between the *Hind*II 8/5 cleavage site and the *Hinf*I 5a/6 cleavage site contains 111 nucleotides. It also contains the *Alu*I cleavage sites between fragments 12b, 15b, 17 and 12a (Fig. 1). The sequence determination allowed the

Present addresses: \*Department of Biochemistry, University of Bristol, Bristol BS8 1TD, UK. †Department of Biochemistry, University of British Columbia, Vancouver BC, Canada V6T 1W5.