

# Sequence of the gene coding for the $\beta$ -subunit of dinitrogenase from the blue-green alga *Anabaena*

(MoFe protein/FeS proteins/cyanobacteria/ribulosebiphosphate carboxylase/endosymbiotic hypothesis)

BARBARA J. MAZUR AND CHOK-FUN CHUI

Central Research and Development Department, E. I. du Pont de Nemours and Company, Wilmington, Delaware 19898

Communicated by H. E. Simmons, Jr., July 23, 1982

**ABSTRACT** The nitrogen fixation *nifK* gene of the blue-green alga *Anabaena*, which codes for the  $\beta$ -subunit of dinitrogenase, has been subjected to sequence analysis. The *nifK* protein is predicted to be 512 amino acids long, to have a  $M_r$  of 57,583, and to contain six cysteine residues. Three of these cysteines are within peptides homologous to FeS cluster-binding cysteinyl peptides from ferredoxins and from a high potential iron protein and, thus, may be ligands to which FeS clusters bind in dinitrogenase. The sequences surrounding the cysteine residues are 70% homologous to the corresponding cysteinyl tryptic peptides of the *Azotobacter vinelandii* dinitrogenase, although the positions of the cysteine residues are not always conserved between the two proteins. A 15-amino acid coding sequence precedes *nifK* on its transcript. Amino acid codon usage is highly asymmetric and parallels that found for the *Anabaena* dinitrogenase reductase gene (*nifH*). Putative promoter and ribosome binding site sequences were identified for *nifK*. These regulatory sequences are homologous to sequences preceding *nifD*; *nifD* codes for the  $\alpha$ -subunit of dinitrogenase but is separated from *nifK* on the chromosome by 11,000 nucleotides. The *nifK* promoter also is virtually identical to a promoter-like sequence that immediately precedes the start of the transcript for the large subunit of ribulosebiphosphate carboxylase from maize chloroplasts. This homology appears to support the theory that chloroplasts evolved from blue-green algae.

Nitrogenase catalyzes the reduction of atmospheric nitrogen to ammonia. The nitrogenase complex is an association between two enzymes—the dinitrogenase, or MoFe protein, and the dinitrogenase reductase, or Fe protein. Dinitrogenase is itself multimeric, being composed of two  $\alpha$ -subunits of  $M_r \approx 56,000$  and two  $\beta$ -subunits of  $M_r \approx 60,000$ . Complexed to it are four  $Fe_4S_4$  clusters and two molecules of a MoFe cofactor; the latter is the active site for nitrogen reduction. Dinitrogenase is reduced by dinitrogenase reductase, a dimer of  $M_r$  33,000 subunits (1).

Seventeen genes are required for nitrogen fixation in *Klebsiella pneumoniae*, the only organism for which a detailed genetic map of the nitrogen fixation (*nif*) genes has been constructed. The *nifK* and *nifD* genes code for the  $\beta$ - and  $\alpha$ -subunits of dinitrogenase, respectively, whereas the *nifH* gene codes for the nitrogenase reductase. The remaining genes are required for activation of these enzymes, for formation of the MoFe cofactor, for electron transport, and for regulation of the *nif* genes (2).

We have been studying the regulation of nitrogen fixation in *Anabaena* 7120, a photosynthetic blue-green alga that can fix simultaneously both carbon and nitrogen. No methods for genetic transformation have been identified for *Anabaena*, hampering construction of a genetic map. This difficulty was

circumvented in the case of the *nif* genes by exploiting a homology between *Anabaena* and *Klebsiella nif* genes. By hybridizing cloned *Anabaena* and *Klebsiella nif* genes to each other, we were able to construct a limited map of the *Anabaena nif* genes (3). Surprisingly, the *Anabaena* and *Klebsiella nif* genes are rearranged relative to each other. Particularly striking are the different locations of the three nitrogenase structural genes, which are linked in a single operon in *Klebsiella* but not in *Anabaena*. Instead, in *Anabaena*, *nifH*—which codes for the dinitrogenase reductase—is linked to *nifD*—which codes for the  $\alpha$ -subunit of dinitrogenase—whereas *nifK*—which codes for the  $\beta$ -subunit of dinitrogenase—is separated from *nifD* by 11 kilobase pairs (kbp) (4).

The *nifK* sequence that we report here reveals one way in which these unlinked genes can be coordinately regulated. The putative promoter sequences for the *nifK* and *nifH*–*nifD* transcripts are homologous, as are their ribosome binding sites. Surprisingly, the putative *nifK* promoter also is virtually identical to a promoter-like sequence that immediately precedes the start of transcription of the gene for the large subunit of ribulosebiphosphate carboxylase (RuP<sub>2</sub>Case) from maize chloroplasts (5), lending support to the theory that chloroplasts evolved from blue-green algae (6).

## MATERIALS AND METHODS

**DNA Sequence Analysis.** Chemical sequence analysis was performed according to the procedures of Maxam and Gilbert (7). Chain termination sequence analysis was performed according to the methods of Sanger *et al.* (8).

## RESULTS AND DISCUSSION

**DNA Sequence of *nifK*.** A map of *nifK* is shown in Fig. 1. The gene was defined through a series of hybridizations between cloned fragments of *Klebsiella nifK* (9) and *Anabaena* DNA. A 1-kbp region of homology was found on three adjacent *Hind*III subfragments of a 17-kbp *Anabaena Eco*RI clone (4). These three *nifK*-containing subfragments, 1, 0.7, and 3 kb in size, were cloned into pBR322 (10) and into M13 mp5 (11) and were subjected to sequence analysis by the chemical and chain termination methods, respectively (7, 8). Overlapping DNA sequences were obtained at all restriction sites, including the *Hind*III sites that separate the subclones. Sequences were confirmed either by sequence analysis of the complementary strand or by using the complementary sequence analysis technique.

The 2,200 base pairs (bp) that flank and include the *nifK* gene are shown in Fig. 2. The direction of transcription of the gene was determined by hybridizing each of the M13 subclones with

The publication costs of this article were defrayed in part by page charge payment. This article must therefore be hereby marked "advertisement" in accordance with 18 U. S. C. §1734 solely to indicate this fact.

Abbreviations: *nif*, nitrogen fixation; kb, kilobase(s); kbp, kilobase pair(s); bp, base pair(s); RuP<sub>2</sub>Case, ribulosebiphosphate carboxylase.



GAGAAAACG CCGCTGGTAG ACGAAAGTGG CTCTAAGTCT GCAAAGGCTT GTCGATATTT GTCTTGACCC TGATTTTGCA TCGCTGTGGT ATTAGCCTAT  
 10 20 30 40 50 60 70 80 90 100  
 ATTTAGCCTA AAAATTAATG TGTTATCAGC AAACAATGTT CATCACTAAC ACTGCTCAGT GCAAACATTA AGCTGTGGAA AGCCATTAAA CCACAAAAAG  
 110 120 130 140 150 160 170 180 190 200  
 GATTACTCCG GCCCTTATCA CGGTTACCAC GGATTGCGTA TCTTCGCCCG TGACATGGAT TTAGCCCTCA ACAGCCCAAC TTGGAGCTTG ATTGGCGCTC  
 210 220 230 240 250 260 270 280 290 300  
 CTTGGAAGAA AGCGGCTGCA AAGGCTAAGG CTGCTGCCTA ATTCCAGGTA AATAGAGGGA TAGCCTGGGG TTGTTGCCCC AGGAACCCAG GGAAGAATGG  
 310 320 330 340 350 360 370 380 390 400  
 TGAATGCTGA ATGCTGGATG AAGATTTTTA TTCATTATTC ATGATTATC ACTCGTTACT ACCTTGAGGG GGAGTGAACC TCCCAGGCTA TCCTCACTCA  
 410 420 430 440 450 460 470 480 490 500  
 TCACTTACAA ACCAACCAGC AAGCGTAGAG AGATACAACA Met Pro Gln Asn Pro Glu Arg Thr Val Asp His Val Asp Leu Phe Lys Gln  
 510 520 530 540 550 560 570 580 590  
 Pro Glu Tyr Thr Glu Leu Phe Glu Asn Lys Arg Lys Asn Phe Glu Gly Ala His Pro Pro Glu Glu Val Glu Arg Val Ser Glu  
 CCA GAA TAC ACC GAG CTA TTT GAA AAC AAG AGA AAG AAC TTT GAA GGC GCT CAT CCT CCT GAA GAA GTT GAA AGA GTG TCT GAA  
 600 610 620 630 640 650 660 670  
 Trp Thr Lys Ser Trp Asp Tyr Arg Glu Lys Asn Phe Ala Arg Glu Ala Leu Thr Val Asn Pro Ala Lys Gly Cys Gln Pro Val  
 TGG ACA AAA TCT TGG GAC TAC CGG GAA AAG AAC TTC GCT CGT GAA GGC GCT TTA ACC GTT AAC CCT GCT AAA GGT TGC CAA CCT GTA  
 680 690 700 710 720 730 740 750  
 Gly Ala Met Phe Ala Ala Leu Gly Phe Glu Gly Thr Leu Pro Phe Val Gln Gly Ser Gln Gly Cys Val Ala Tyr Phe Arg Thr  
 GGC GCG ATG TTC GCT GCT TTG GGT TTT GAA GGT ACT CTA CCT TTC GTA CAA GGT TCT CAA GGT TGC GTT GCT TAC TTC CGT ACA  
 770 780 790 800 810 820 830 840  
 His Leu Ser Arg His Tyr Lys Glu Pro Cys Ser Ala Val Ser Ser Ser Met Thr Glu Asp Ala Ala Val Phe Gly Gly Leu Asn  
 CAC CTC AGC CGT CAC TAC AAA GAG CCT TGC TCC GCA GTA TCT TCT TCC ATG ACA GAA GAT GCA GCA GTA TTC GGT GGT TTG AAC  
 850 860 870 880 890 900 910 920  
 Asn Met Ile Glu Gly Met Gln Val Ser Tyr Gln Leu Tyr Lys Pro Lys Met Ile Ala Val Cys Thr Cys Met Ala Glu Val  
 AAC ATG ATC GAA GGT ATG CAG GTT TCA TAC CAA CTG TAC AAG CCT AAG ATG ATT GCT GTT TGC ACC ACC TGT ATG GCG GAA GTT  
 930 940 950 960 970 980 990 1000 1010  
 Ile Gly Asp Asp Leu Gly Ala Phe Ile Thr Asn Ser Lys Asn Ala Gly Ser Ile Pro Gln Asp Phe Pro Val Pro Phe Ala His  
 ATC GGA GAT GAC TTG GGC GCG TTC ATC ACC AAC TCC AAG AAC GCT GGT TCT ATT CCT CAA GAT TTC CCC GTA CCC TTT GCT CAC  
 1020 1030 1040 1050 1060 1070 1080 1090  
 Thr Pro Ser Phe Val Gly Ser His Ile Thr Gly Tyr Asp Asn Met Lys Gly Ile Leu Ser Asn Leu Thr Glu Gly Lys Lys  
 ACA CCT AGC TTT GTT GGT TCC CAC ATC ACT GGC TAC GAC AAC ATG AAG GAT ATT CTG TCT AAC TTG ACA GAA GGT AAG AAC  
 1100 1110 1120 1130 1140 1150 1160 1170  
 Lys Ala Thr Ser Asn Gly Lys Ile Asn Phe Ile Pro Gly Phe Asp Thr Tyr Val Gly Asn Asn Arg Glu Leu Lys Arg Met Met  
 AAA GCT ACC AGC AAC GGC AAA ATT AAC TTC ATT CCT GGT TTT GAT ACC TAT GTA GGT AAC AAC CGC GAA TTG AAG CGG ATG ATG  
 1190 1200 1210 1220 1230 1240 1250 1260  
 Gly Val Met Gly Val Asp Tyr Thr Ile Leu Ser Asp Ser Ser Asp Tyr Phe Asp Ser Pro Asn Met Gly Glu Tyr Glu Met Tyr  
 GGT GTA ATG GGT GTT GAC TAC ACC ATC CTG TCT GAC AGC AGC GAC TAT TTT GAT TCA CCT AAC ATG GGT GAA TAC GAA ATG TAC  
 1270 1280 1290 1300 1310 1320 1330 1340  
 Pro Ser Gly Thr Lys Leu Glu Asp Ala Ala Asp Ser Ile Asn Ala Lys Ala Thr Val Ala Leu Gln Ala Tyr Thr Thr Pro Lys  
 CCA AGT GGT ACA AAG CTG GAA GAT GCG GCT GAT TCT ATC AAC GCT AAA GCA ACT GTT GCT CTC CAA GCT TAC ACC ACA CCT AAG  
 1360 1370 1380 1390 1400 1410 1420 1430  
 Thr Arg Glu Tyr Ile Lys Thr Gln Trp Lys Gln Glu Thr Gln Val Leu Arg Pro Phe Gly Val Lys Gly Thr Asp Glu Phe Leu  
 ACC CGC GAA TAC ATC AAA ACC CAG TGG AAG CAA GAA ACA GTA TTG CGC CCC TTC GGT ATT AAG GGT ACT GAC GAG TTC TTG  
 1440 1450 1460 1470 1480 1490 1500 1510  
 Thr Ala Val Ser Glu Leu Thr Gly Lys Ala Ile Pro Glu Glu Leu Glu Ile Glu Arg Gly Arg Leu Val Asp Ala Ile Thr Asp  
 ACT GCT GTT TCT GAA TTG ACC GGT AAA GCT ATT CCT GAA GAA TTG GAA ATC GAA CGC GGT CGT TTA GTT GAT GCT ATC ACT GAC  
 1520 1530 1540 1550 1560 1570 1580 1590  
 Ser Tyr Ala Trp Ile His Gly Lys Lys Phe Ala Ile Tyr Gly Asp Pro Asp Leu Ile Ile Ser Ile Thr Ser Phe Leu Leu Glu  
 TCC TAC GCT TGG ATT CAT GGT AAG AAG TTC GCT ATC TAC GGC GAT CCA ATC TTG ATC ATC TCC ATC ACC AGC TTC TTG TTA GAA  
 1610 1620 1630 1640 1650 1660 1670 1680  
 Met Gly Ala Glu Pro Val His Ile Leu Cys Asn Asn Gly Asp Asp Thr Phe Lys Lys Glu Met Glu Ala Ile Leu Ala Ala Ser  
 ATG GGT GCT GAA CCA GTA CAC ATC CTC TGC AAC AAC GGT GAT GAC ACC TTC AAG AAA GAA ATG GAA GCT ATC CTC GCT GCT AGC  
 1690 1700 1710 1720 1730 1740 1750 1760  
 Pro Phe Gly Lys Glu Ala Lys Val Trp Ile Gln Lys Asp Leu Trp His Phe Arg Ser Leu Leu Phe Thr Glu Pro Val Asp Phe  
 CCA TTT GGT AAA GAA GCT AAA GTC TGG ATT CAA AAA GAC TTG TGG CAC TTC CGT TCC TTG TTG TTC ACC GAG CCI GTA GAC TTC  
 1780 1790 1800 1810 1820 1830 1840 1850  
 Phe Ile Gly Asn Ser Tyr Gly Lys Tyr Leu Trp Arg Asp Thr Ser Ile Pro Met Val Arg Ile Gly Tyr Pro Leu Phe Asp Arg  
 TTC ATC GGT AAC TCC TAC GGT AAG TAC CTG TGG CGC GAT ACC AGC ATC CCA ATG GTG CGG ATT GGT TAT CCT CTC TTC GAT CGC  
 1860 1870 1880 1890 1900 1910 1920 1930  
 His His Leu His Arg Tyr Ser Thr Leu Gly Tyr Gln Gly Gly Asn Ile Leu Asn Trp Val Val Asn Thr Leu Leu Asp Glu  
 CAC CAC TTA CAC CGC TAT TCT ACC CTC GGC TAC CAA GGT GGT CTA AAT ATC CTC AAC TGG GTT GGT AAC ACC CTG TTG TAT GAA  
 1940 1950 1960 1970 1980 1990 2000 2010  
 Met Asp Arg Ser Thr Asn Ile Thr Gly Lys Thr Asp Ile Ser Phe Asp Leu Ile Arg  
 ATG GAT CGC AGC ACC AAC ATC ACT GGT AAG ACC GAT ATC TCC TTT GAC TTG ATC CGC TAGA AATTAATGCA GCGTGCCATT  
 2030 2040 2050 2060 2070 2080 2090 2100  
 GAAAGGTAGA ACTTAGGGAC TGGGGATTGG GTATTGGGTA CTAGGAATAT TATCTTCCCA GTCCCTTCCA GTCCCCGAGA CCCTTTG  
 2110 2120 2130 2140 2150 2160 2170 2180

FIG. 2. The nucleotide sequence of *nifK* of *Anabaena*. The derived amino acid sequence of the  $\beta$ -subunit of dinitrogenase is also shown above residues 323-1,861.



and two more are at different positions. One of the missing cysteines is within a short peptide that is not homologous to the *Anabaena* protein, and the other is within a homologous region of the *Anabaena* protein but is an alanine in *Anabaena* (residue 384, nucleotide 1,690). The *Anabaena* cysteine at residue 111 (nucleotide 871) is a valine in *Azotobacter*, whereas *Azotobacter*, instead, has a cysteine two residues downstream, again at the position of an *Anabaena* alanine. Likewise, the first cysteine in the peptide Cys-Thr-Thr-Cys (Cys-150, nucleotide 988) is missing in *Azotobacter*, whereas *Azotobacter* once again has a cysteine at the position of an *Anabaena* alanine (residue 78, nucleotide 772).

The absence of this Cys-X-X-Cys peptide in *Azotobacter* is somewhat surprising, in that the cysteines of such peptides ligand FeS clusters in ferredoxins and in a high potential iron protein (24, 25) and also might be expected to ligand FeS clusters in dinitrogenase. One other *Anabaena* tetrapeptide surrounding a cysteine also is homologous to an FeS cluster binding cysteine in a high potential iron protein, the peptide Lys-Gly-Cys-Gln (Cys-70, nucleotide 748). The cysteine in this peptide is conserved in *Azotobacter*, although the adjacent glycine is instead an alanine. The adjacent glutamine is followed by Pro-Val; Lundell and Howard have postulated that such Pro-Val residues should provide an important secondary structure for ligand folding around FeS centers (23). It will be of interest to learn whether either of the translocated cysteines can perform equivalent functions in the *Anabaena* and *Azotobacter* proteins and to learn which cysteines in fact ligand FeS clusters. Unlike the cysteinyl peptides, the NH<sub>2</sub>-more and COOH-terminal peptides are not homologous between the *Anabaena* and *Azotobacter*  $\beta$ -subunit proteins (20).

**Codon Usage in the *nifK* Gene.** The codon usage for the *nifK* protein is strikingly asymmetric. Fig. 4 tabulates the usage and compares it to the only other blue-green algal protein that has been subjected to sequence analysis, *nifH* (12). There are nine codons that are never used in *nifK*, and five of these also are never used in *nifH*. An additional eight codons are used but once or twice in *nifK* and *nifH*. Many of the remaining codons are used only a few times, or else heavily. This bias is not peculiar to the *nif* genes, as *nifH* from *Klebsiella* exhibits a different pattern of codon usage (17). Instead, it may reflect the distribution of tRNA populations in algal cells or in their heterocysts (26).

We thank Steven Robinson and Robert Haselkorn for permission to cite their unpublished data.

1. Mortenson, L. E. & Thorneley, R. N. F. (1979) *Annu. Rev. Biochem.* **48**, 387-418.
2. Roberts, G. P. & Brill, W. J. (1981) *Annu. Rev. Microbiol.* **35**, 207-235.
3. Mazur, B. J., Rice, D. & Haselkorn, R. (1980) *Proc. Natl. Acad. Sci. USA* **77**, 186-190.
4. Rice, D., Mazur, B. J. & Haselkorn, R. (1982) *J. Biol. Chem.* **257**, in press.
5. McIntosh, L., Poulsen, C. & Bogorad, L. (1980) *Nature (London)* **288**, 556-560.
6. Ris, H. & Plaut, W. (1962) *J. Cell Biol.* **13**, 383-391.
7. Maxam, A. M. & Gilbert, W. (1980) *Methods Enzymol.* **65**, 499-560.
8. Sanger, F., Nicklen, S. & Coulson, A. R. (1977) *Proc. Natl. Acad. Sci. USA* **74**, 5463-5467.
9. Riedel, G. E., Ausubel, F. M. & Cannon, F. C. (1979) *Proc. Natl. Acad. Sci. USA* **76**, 2866-2870.
10. Bolivar, F., Rodriguez, R., Greene, P., Betlach, M., Heyneker, H., Boyer, H., Crosa, J. & Falkow, S. (1977) *Gene* **2**, 95-113.
11. Messing, J. (1979) *Recomb. DNA Tech. Bull.* **2**, 43-48.
12. Mevarech, M., Rice, D. & Haselkorn, R. (1980) *Proc. Natl. Acad. Sci. USA* **77**, 6476-6480.
13. Rosenberg, M. & Court, D. (1979) *Annu. Rev. Genet.* **13**, 319-353.
14. Zurawski, G., Perrot, B., Bottomley, W. & Whitfield, P. R. (1981) *Nucleic Acids Res.* **9**, 3251-3270.
15. Shine, J. & Dalgarno, L. (1974) *Proc. Natl. Acad. Sci. USA* **71**, 1342-1346.
16. Tinoco, I., Jr., Borer, P., Dengler, B., Levine, M. D., Uhlenbeck, O. C., Crothers, D. M. & Gralla, J. (1973) *Nature (London)* **246**, 40-41.
17. Scott, K. F., Rolfe, B. G. & Shine, J. (1981) *J. Mol. Appl. Genet.* **1**, 71-81.
18. Fleming, H. & Haselkorn, R. (1974) *Cell* **3**, 159-170.
19. Kennedy, C., Eady, R. R., Kondorosi, E. & Rekosch, D. K. (1976) *Biochem. J.* **155**, 383-389.
20. Lundell, D. J. & Howard, J. B. (1978) *J. Biol. Chem.* **253**, 3422-3426.
21. McMahon, J. E. & Tinoco, I., Jr. (1978) *Nature (London)* **271**, 275-277.
22. Segrest, J. P. & Feldmann, R. J. (1974) *J. Mol. Biol.* **87**, 853-858.
23. Lundell, D. J. & Howard, J. B. (1981) *J. Biol. Chem.* **256**, 6385-6391.
24. Adman, E. T., Sieker, L. C. & Jensen, L. H. (1973) *J. Biol. Chem.* **248**, 3987-3996.
25. Carter, C. W., Kraut, J., Freer, S. T., Xuong, N.-H., Alden, R. A. & Bartsch, R. G. (1974) *J. Biol. Chem.* **249**, 4212-4225.
26. Ikemura, T. (1981) *J. Mol. Biol.* **146**, 1-21.