# Molecular Characterization of the *Clostridium difficile* Toxin A Gene

C. H. DOVE, S.-Z. WANG, S. B. PRICE,† C. J. PHELPS, D. M. LYERLY, T. D. WILKINS, AND J. L. JOHNSON*

*Department of Anaerobic Microbiology, Virginia Polytechnic Institute and State University, Blacksburg, Virginia 24061*

The gene encoding the toxin A protein of *Clostridium difficile* (strain VPI 10463) was cloned and sequenced. The coding region of 8,133 base pairs has a mol% G+C of 26.9 and encodes 2,710 amino acids. The deduced polypeptide has a molecular mass of ca. 308 kilodaltons. Nearly a third of the gene, at the 3' end, consists of 38 repeating sequences. The repeating units were grouped into two classes, I and II, on the basis of length and the low levels of DNA sequence similarities between them. There were seven class I repeating units, each containing 90 nucleotides, and 31 class II units, which, with two exceptions, were either 60 or 63 nucleotides in length. On the basis of DNA sequence similarities, the class II repeating units were further segregated into subclasses: 7 class IIA, 13 class IIB, 5 class IIC, and 6 class IID. The dipeptide tyrosine-phenylalanine was found in all 38 repeating units, and other amino acid sequences were unique to a specific class or subclass. This region of the protein has epitopes for the monoclonal antibody PCG-4 and includes the binding region for the Galα1-3Galβ1-4GlcNAc carbohydrate receptor. Located 1,350 base pairs upstream from the toxin A translation start site is the 3' end of the toxin B gene. Between the two toxin genes is a small open reading frame, which encodes a deduced polypeptide of ca. 16 or 19 kilodaltons. The role of this open reading frame is unknown.

*Clostridium difficile* is the major causative agent of pseudomembranous colitis in humans (2). The organism produces two toxins, designated toxin A and toxin B (1, 32, 33). They are both cytotoxic and lethal for animals, although toxin B is about 1,000-fold more cytotoxic than toxin A for most cell lines. Both toxins appear to be produced in all toxigenic strains; however, the toxicity of strains may vary by several orders of magnitude (20, 32). The actions of these toxins appear to be quite complex and at present are not understood. Although toxin A has a direct toxic effect on the intestinal mucosa, toxin B does not cause a significant response when given intragastrically to hamsters, unless it is initially mixed with a small amount of toxin A (19). Alternatively, toxin B is also toxic if it is given to hamsters with bruised (injured) ceca. The results are consistent with the initial binding and primary tissue damage being caused by toxin A or by mechanical injury, followed by the entry of toxin B.

Investigators in a number of laboratories have worked on the isolation and physical properties of the toxins (1, 2, 27, 30, 32). Both toxins have been purified to homogeneity (18, 32). An interesting and controversial property of the toxins has been their molecular weights. Initial molecular weight estimations obtained by using native proteins have ranged from 440,000 to 600,000 for toxin A and 360,000 to 500,000 for toxin B (1, 30, 32). However, in later studies there has been controversy as to whether the toxins dissociate into smaller subunits under denaturing conditions. Under these conditions, size estimations for the toxins range from 300,000 (18) to 50,000 (27, 30) and down to 42,500 and 16,000 (29). These discrepancies have been summarized by Lyerly et al. (16) and are difficult to explain (35). Perhaps in some of the isolation procedures a smaller contaminating protein copurified with the toxins and tended to mask the toxins in the polyacrylamide gels. The best approach to resolve this controversy is to clone and sequence the toxin genes.

Several investigators have begun cloning these genes. Muldrow and his collaborators (26) have reported the cloning of a 0.3-kilobase-pair (kb) fragment of the toxin A gene in the lambda bacteriophage expression vector gt11. The expressed peptides reacted with toxin A polyvalent antisera. When the cloned fragment from toxin A was used as a labeled probe, it reacted with a *Pst*I-generated fragment of *C. difficile* DNA, which they estimated as 4.5 kb. We have cloned a 4.7-kb *Pst*I fragment into a plasmid vector (28). This fragment has an internal *Pst*I site which is protected from digestion in the *C. difficile* DNA. When this fragment is expressed, the peptide reacts with both toxin A affinity-purified polyclonal antisera and with the monoclonal antibody PCG-4 (17). Preliminary results on the sequencing of this fragment have shown that there are many repeating sequence units within the fragment (14). Eichel-Streiber et al. (8) have recently cloned portions of the 4.7-kb *Pst*I fragment into a plasmid expression vector and obtained an expression product that also reacted with toxin A antisera. Wren et al. (36) have reported the cloning of toxin A in lambda phage. The clone expressed a protein that caused elongation of Chinese hamster ovary cells, and this protein had an estimated molecular weight of 235,000.

In this study, we have completed the cloning and sequencing of toxin A and its flanking regions. Nearly one-third of the gene (from the 3' end) consists of a series of repeating units which appear to code for the receptor portion of the toxin.

## MATERIALS AND METHODS

**Bacteria, bacteriophages, and plasmids.** DNA isolated from *C. difficile* VPI 10463 was used for cloning. Plasmids pBR322, pUC18, and pUC19 were used for the primary cloning of *C. difficile* DNA fragments. Subclones in the M13 phages mp18 and mp19 were used for DNA sequencing. The plasmids, phages, and *Escherichia coli* host strains JM109, DH5α, and DH5αF' were all obtained from Bethesda Research Laboratories, Inc. *E. coli* strain Chi 1776 was purchased from the American Type Culture Collection.

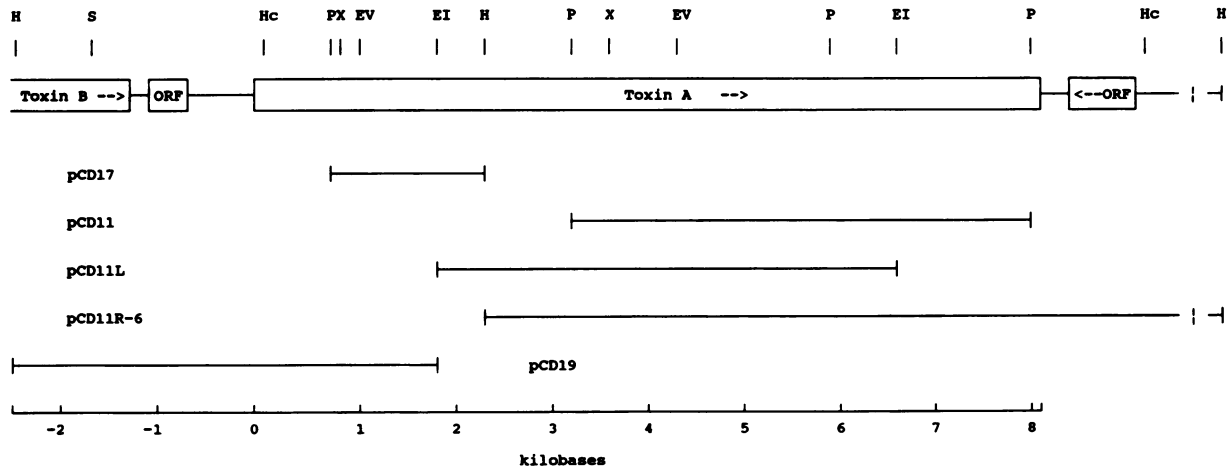**Enzymes and radiolabeled compounds.** Restriction endonu-

---

FIG. 1. Partial endonuclease restriction map of the cloned toxin A region from *C. difficile* strain 10463. Also shown are the sizes and locations of primary clones pCD11, pCD11L, pCD11R-6, pCD17, and pCD19.

clease enzymes were purchased from Bethesda Research Laboratories, Inc., International Biotechnologies, Inc., or Promega Biotec. Endonucleases III and VII, *E. coli* DNA polymerase I, and the polymerase I Klenow fragment were obtained from Bethesda Research Laboratories, Inc. DNA ligase and calf alkaline phosphatase were purchased from Boehringer Mannheim Biochemicals. The T7 DNA polymerase, Sequenase, was purchased from U.S. Biochemical. The labeled nucleotide triphosphate [α-$^{35}$S]dATP was obtained from New England Nuclear. Random primer labeling kits were obtained from International Biotechnologies, Inc. All of the enzymes were used according to the instructions provided by the manufacturers.

**DNA isolations.** High-molecular-weight *C. difficile* DNA was isolated by using a variation of the Marmur procedure (13, 22). The harvested cells (12) were suspended in 50 mM Tris–1 mM EDTA buffer (pH 8.0). After 50 μg of lysozyme per ml was added, the cell suspension was incubated at 37°C until the cells were susceptible to sodium dodecyl sulfate disruption. At the time of disruption, EDTA was added to a final concentration of 40 mM, proteinase K was added to a final concentration of 30 to 50 μg/ml, and β-mercaptoethanol was added to a concentration of 1% to inhibit endogenous nuclease activity in the lysate. Plasmid DNA and the replicating-form DNA of M13 phage were isolated by the Birnboim and Doly alkaline lysis procedure (4). DNA preparations used for probe fragment generation and nested deletions were further purified by CsCl centrifugation. Specific DNA restriction fragments were separated from others on low-melting-point agarose, and the individual bands were cut from the gel for use in the random priming labeling procedure. When restriction fragments of *C. difficile* DNA were needed in a particular size range (i.e., for cloning by chromosome walking), the gel was cut at the lower size range and this part was removed. A well was then cut into the gel at the upper size range, the polarity of the electrophoresis unit was reversed, and the fragments were electroeluted into the well. This tended to concentrate the fragments entering the well and resulted in a lower eluate volume. The fragments were then ethanol precipitated.

**Primary cloning.** Of the four primary clones used in this study (Fig. 1), three were cloned by the chromosome walking approach. The 2.6-kb *Pst*I fragment of pCD11 was used as a probe for cloning pCD11L, the 0.5-kb *Hin*dIII-*Eco*RI fragment of pCD11L was used as a probe for cloning pCD17,

which was then used as a probe to clone pCD19, and pCD11 was used to detect clone pCD11R-6. The cloning was carried out under EK-2, BL-2 containment with *E. coli* Chi 1776 as host. Cells were made competent by the Hanahan procedure (11, 21).

**Toxicity assays.** Lysates from each primary clone were checked for animal and cell toxicity. Mouse lethality tests were performed by injecting five 8-week-old BALB/c mice (Dominion Laboratories, Dublin, Va.) intraperitoneally with 200 μl of lysate and observing them for illness or death. Cytotoxicity was checked in the Chinese hamster ovary (CHO) cell assay by following a procedure previously described (7). Lysates of *E. coli* Chi 1776 transformed with pUC18 were used as negative controls.

**DNA sequencing and sequence analysis.** Both strands of the DNA were sequenced by using the dideoxy-chain termination procedure developed by Sanger et al. (3, 31). DNA fragments were cloned into M13, and nested deletions were generated in replicating-form DNA by using the exonuclease III and exonuclease VII procedure (37). Restriction sites used for subcloning were sequenced across by using oligonucleotide primers and double-stranded sequencing. Synthetic oligonucleotide primers were also used for filling occasional sequence gaps not covered by the nested deletions.

Sequence analysis was done by using the Pustell programs from International Biotechnologies, Inc., and the Sequence Analysis Software Package from the Genetics Computer Group, University of Wisconsin. The data bases that were searched included the GenBank data base and the National Biomedical Research Foundation Protein Sequence Data Base. Unweighted Pair Group cluster analysis was done by using the NTSYS-pc programs (F. J. Rohlf, Exeter Publishing, Ltd.).

**N-terminal sequencing.** Toxin A was purified from culture filtrates of *C. difficile* VPI 10463 by sequential ammonium sulfate precipitation, ion-exchange chromatography, and precipitation at pH 5.6 as previously described (18). The highly purified protein was denatured with a final concentration of 2.5% sodium dodecyl sulfate–5% 2-mercaptoethanol at 100°C for 2 min and subjected to sodium dodecyl sulfate-polyacrylamide gel electrophoresis. After electrophoresis, the protein was transferred to polyvinylidene difluoride membranes by electroblotting and the N-terminal amino acid

FIG. 2. Nucleotide and deduced amino acid sequences of *C. difficile* toxin A gene.

FIG. 2—*Continued*

FIG. 3. Nucleotide sequence similarity cluster analysis of the class I repeating sequences.

sequence was determined by previously described methods (23).

## RESULTS

**Primary clones of toxin A.** Relationships between the five primary clones, each containing a portion of the *C. difficile* toxin A gene, are shown in Fig. 1. Also included in the figure is a partial restriction map of this 15-kb region of the *C. difficile* genome. Clone pCD11 has been partially characterized and shown to contain a carbohydrate binding region and antigenic epitopes which react with the monoclonal antibody PCG-4 (28). Clone pCD11R-6, in addition to containing the entire pCD11 insert and most of the pCD11L insert, contains the last 80 bases of the toxin A gene and approximately 4.1 kb of additional sequences downstream from the toxin A gene. The downstream region contains two open reading frames (ORFs) and part of the third, one of which is shown in Fig. 1. All of these ORFs read in the direction opposite that of the toxin A gene (data not shown). Clone pCD11L contains an additional 1.5 kb of sequence upstream of the pCD11 insert. Clone pCD17 was used as a probe for cloning pCD19. Clone pCD19 codes for the 5' end of toxin A, a small ORF that could code for a 16- or 19-kilodalton (kDa) protein and 1.2 kb of toxin B. These clones were not toxic for mice or CHO cells. The clone immediately upstream from the pCD19 insert was found to contain the remainder of the toxin B gene, and we have since been able to reconstruct the intact gene in a plasmid. The recombinant protein expressed by this plasmid is cytotoxic to tissue cells, is lethal to mice, and has immunological identity with toxin B (D. M. Lyerly and J. L. Johnson, unpublished data).

**Nucleotide and amino acid sequences for toxin A.** The nucleotide sequence and the deduced amino acids for the toxin A gene are shown in Fig. 2 (GenBank accession number, M30307). The open reading frame is 8,133 nucleotides long and codes for 2,710 amino acids. The gene contains 26.9 mol% G+C, and the deduced protein has a molecular mass of 308,103 Da. The amino acid sequence of the N-terminal end of toxin A was determined by microanalysis after electrophoresis under denaturing conditions (23), and the first 10 amino acids agree with the first 10 deduced amino acids of the toxin A open reading frame, indicating that there are no posttranslational modifications involving a signal peptide.

An interesting property of this gene is the repeating sequences at the 3' end. A total of 2,551 nucleotides, or 31.5% of the gene, are in 38 contiguous repeating units. This region extends from nucleotides 5,545 to 8,106. The repeating units were grouped into two classes, I and II, on the basis of the low levels of DNA sequence similarities between



FIG. 4. Nucleotide sequence similarity cluster analysis of the class II repeating sequences.

them. There are 7 class I and 31 class II repeating units. Each of the class I repeats is 90 nucleotides long, and the class II repeats are either 60 or 63 nucleotides long, with the one exception being 66 nucleotides long. The class II repeats have been subdivided into 7 class IIA, 13 class IIB, 5 class IIC, and 6 class IID repeats.

Nucleotide sequence similarities among the class I repeats are shown in Fig. 3. Similarities ranged from 73 to 98%, with the average values in the cluster analysis being 80% or greater. Nucleotide sequence similarities among the class II repeats are shown in Fig. 4. With the exception of class IID, clustering within each subclass is high, being 70% or higher for class IIA, 65% or higher for class IIB, and 76% or higher for class IIC repeats. The class IID repeats are a diverse collection, in that all are very distinct. Two of them fit closer to the class IIB cluster and one fits closer to the class IIC group than to the others in class IID. This is also the only group in which there is any size variation; repeat unit class $IID_4$ has an extra AAA codon, while $IID_5$ has one fewer codon.

The deduced amino acid residues for the repeated se-

**Class I peptides**

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| I₁ | 1891-1920 | M Q L | G V F | K G P D | G F E Y F A P A N T | Q N | N N I E G Q A I |
| I₂ | 2025-2054 | V K I | G V F | S T S N | G F E Y F A P A N T | Y N | N N I E G Q A I |
| I₃ | 2159-2188 | M Q I | G V F | K G P N | G F E Y F A P A N T | D A | N N I E G Q A I |
| I₄ | 2273-2302 | M V T | G V F | K G P N | G F E Y F A P A N T | H N | N N I E G Q A I |
| I₅ | 2407-2436 | M Q I | G V F | K G P N | G F E Y F A P A N T | D A | N N I E G Q A I |
| I₆ | 2520-2549 | M Q I | G V F | K G P D | G F E Y F A P A N T | D A | N N I E G Q A I |
| I₇ | 2611-2640 | P Q I | G V F | K G S N | G F E Y F A P A N T | D A | N N I E G Q A I |

**CONSENSUS**    M Q I G V F K G P N G F E Y F A P A N T D A N N I E G Q A I

FIG. 5. Deduced amino acid sequences for the class I repeating units. Unit designations (I₁ to I₇) are listed in order from the N-terminal to C-terminal direction. The inclusive amino acid residue numbers are given for each unit, and the conserved amino acids are boxed.

**CLASS IIA PEPTIDES**

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| A₁ | 1921-1940 | V | Y Q | S K F L | T | L | N G K K | Y Y F | D N N | S |
| A₂ | 2055-2074 | V | Y Q | S K F L | T | L | N G K K | Y Y F | D N N | S |
| A₃ | 2189-2208 | L | Y Q | N E F L | T | L | N G K K | Y Y F | G S D | S |
| A₄ | 2303-2322 | V | Y Q | N K F L | T | L | N G K K | Y Y F | D N D | S |
| A₅ | 2437-2456 | L | Y Q | N K F L | T | L | N G K K | Y Y F | G S D | S |
| A₆ | 2550-2569 | R | Y Q | N R F L | Y | L | H D N I | Y Y F | G N N | S |
| A₇ | 2641-2660 | R | Y Q | N R F L | H | L | L G K I | Y Y F | G N N | S |

**CONSENSUS**   V Y Q N K F L T L N G K K Y Y F G N N S

**CLASS IIB PEPTIDES**

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| B₁ | 1849-1869 | N L V | T G | W Q T | I | N G K K | Y Y F | D I N T G |
| B₂ | 1941-1961 | K A V | T G | W R I | I | N N E K | Y Y F | N P N N A |
| B₃ | 2075-2095 | K A V | T G | W Q T | I | D S K K | Y Y F | N T N T A |
| B₄ | 2096-2116 | E A A | T G | W Q T | I | D G K K | Y Y F | N T N T A |
| B₅ | 2117-2137 | E A A | T G | W Q T | I | D G K K | Y Y F | N T N T A |
| B₆ | 2209-2229 | K A V | T G | W R I | I | N N K K | Y Y F | N P N N A |
| B₇ | 2323-2343 | K A V | T G | W Q T | I | D G K K | Y Y F | N L N T A |
| B₈ | 2344-2364 | E A A | T G | W Q T | I | D G K K | Y Y F | N L N T A |
| B₉ | 2365-2385 | E A A | T G | W Q T | I | D G K K | Y Y F | N T N T F |
| B₁₀ | 2457-2477 | K A V | T G | L R T | I | D G K K | Y Y F | N T N T A |
| B₁₁ | 2478-2498 | V A V | T G | W Q T | I | N G K K | Y Y F | N T N T S |
| B₁₂ | 2570-2590 | K A A | T G | W V T | I | D G N R | Y Y F | E P N T A |
| B₁₃ | 2661-2681 | K A V | T G | W Q T | I | N G K V | Y Y F | M P D T A |

**CONSENSUS**   K A V T G W Q T I D G K K Y Y F N T N T A

**CLASS IIC**

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| C₁ | 1870-1890 | A | A | L T S | Y | K I | I | N | G K H F Y F | N N D | G V |
| C₂ | 2004-2024 | I | A | F N G | Y | K T | I | D | G K H F Y F | D S D | C V |
| C₃ | 2138-2158 | I | A | S T G | Y | T I | I | N | G K H F Y F | N T D | G I |
| C₄ | 2386-2406 | I | A | S T G | Y | T S | I | N | G K H F Y F | N T D | G I |
| C₅ | 2499-2519 | I | A | S T G | Y | T I | I | S | G K H F Y F | N T D | G I |

**CONSENSUS**   I A S T G Y T I I N G K H F Y F N T D G I

**CLASS IID**

| | | | | | |
|---|---|---|---|---|---|
| D₁ | 1962-1982 | I A A V G L Q V I D N N K Y | Y F | N P D T A |
| D₂ | 1983-2003 | I I S K G W Q T V N G S R Y | Y F | D T D T A |
| D₃ | 2230-2250 | I A A I H L C T I N N D K Y | Y F | S Y D G I |
| D₄ | 2251-2271 | L Q N G Y I T I E R N N F | Y F | D A N N E S K |
| D₅ | 2591-2611 | M G A N G Y K T I D N K N F | Y F | R N G L |
| D₆ | 2682-2702 | M A A A G G L F E I D G V I | Y F | F G V D G |

**CONSENSUS**   I A A - G - - T I - N - - Y Y F - - - D - -

quences are shown in Fig. 5 and 6. The inclusive amino acid residue numbers are given for each repeat unit, and the conserved amino acid residues within each class or class subgroup are boxed. Seventy percent of the amino acids in the class I peptides are conserved among the units, while less than 50% are conserved within each of the class II subgroups. The dipeptide tyrosine-phenylalanine (YF) is the most conserved and can be found in all 38 repeat units. It represents residues 14 and 15 in the class I units and residues 15 and 16 in the class II repeats, except for unit IID₄. Base differences in the regions of conserved amino acids involved the codon's third base as expected, whereas switching from one amino acid to another in a given position usually involved a total codon change or at least two of the bases.

A hydropathic index plot for the deduced toxin A protein and a map of the repeat units are shown in Fig. 7. There is no evidence for a signal peptide at the amino-terminal end; this finding is in agreement with the lack of posttranslational modification of the N-terminal end of the protein. The only strongly hydrophobic region in the deduced protein is from residues 1,050 to 1,100. There appears to be a periodicity in the hydropathic index within each repeat region. However, the repeat region is for the most part hydrophilic.

The 160 bases immediately upstream from the toxin A translation initiation site are shown in Fig. 2. There appears to be a ribosomal binding site (GGAGGT) starting six bases upstream of the initiation codon. Since we do not know where transcription initiates, it is difficult to predict promoter regions, although there are several TA-rich areas in the region of 160 bases upstream (Fig. 2). Other than these, there do not appear to be any other unique structures, such as inverted or tandem repeats.

**Small protein.** A small ORF (ca. 500 base pairs; Fig. 8; GenBank accession number, 30308) is located 122 bases downstream from the stop codon of the toxin B gene. Although the deduced amino acid sequence begins with the first start codon, there are two additional ATG codons at the amino acid residue positions 25 and 27. There appear to be ribosomal binding sites in the −10 regions of the first (GGTGGA) and third (GGAGGC) ATG codons. The deduced protein would have a molecular mass of 18,798 Da by using the longer sequence and a 15,878-Da molecular mass by using the shorter sequence. The pI values for the two peptides are 9.22 and 9.11, respectively. The hydropathic

FIG. 6. Deduced amino acid sequences for the class II repeating units. Unit designations are made in the same manner as for the class I units.

FIG. 7. Hydropathy plot and repeating unit map for *C. difficile* toxin A gene. Hydrophobic regions are indicated by positive values.

indexes were determined for both versions of the ORF (data not included). The deduced peptide is in general hydrophilic, and there does not appear to be a signal peptide in the first 25 amino acid residues; however, for a polypeptide starting at amino acid residue 27 (the third ATG codon), there is a short hydrophobic region that is characteristic of other signal sequences (24).

## DISCUSSION

We report here the molecular mass of 308,103 Da for the deduced toxin A protein of *C. difficile*. This is in agreement with previous studies that reported a large size for this toxin (1, 18, 32, 33). Although we have not been able to express toxicity from the cloned fragments, the 2.1-kb *Pst*I fragment at the 3′ end of the gene has been used to express the major antigenic and carbohydrate binding sites of the toxin (15, 28). In fact, antiserum against this portion of the protein neutralizes the enterotoxicity of toxin A, and this is further evidence that the repeating units represent the binding portion (D. M. Lyerly and T. D. Wilkins, unpublished data).

The mechanism of action of toxin A is unknown. In the data base searches, we were unable to find any amino acid sequence similarities with other characterized toxins or enzymes. We cannot rule out a second peptide associating with this one, for example, the small ORF protein. The loss of such a protein would have very little effect on the electrophoretic migration and probably would not be detected if the protein existed in equimolar amounts with the large protein. Also, after electrophoresis under denaturing conditions, only antigenicity has been measured and not toxicity. It is not yet known whether the small protein is even expressed in *C. difficile*, so any presumed role for the small protein in the toxicity of the organism will have to await further study.

The most interesting feature of the toxin A gene is the repeating sequences in the carbohydrate binding region, which, as seen by the hydropathy plot, contains the most hydrophilic portion of the molecule. This is at the carboxyl end of the protein and includes over a third of the polypeptide. Proteins with repeating units have been reported from a wide range of organisms. Some of the highly antigenic

```
ATAAAAATAT GTTAAATATA TCCTCTTATA CTTAAATATA TAAAAATAAA CAAAATGATA        60

CACTACATAA AGTGTTCTAT CTAATATGAA GATTTACCAA TAAAAAGGTG GACTATGATG       120

A ATG CAC AGT AGT TCA CCT TTT TAT ATT TCT AAT GGT AAC AAA ATA TTT TTT   172
  M   H   S   S   S   P   F   Y   I   S   N   G   N   K   I   F   F

TAT ATA AAC CTA GGA GGC GTT ATG AAT ATG ACA ATA TCT TTT TTA TCA GAG     223
 Y   I   N   L   G   G   V   M   N   M   T   I   S   F   L   S   E

CAT ATA TTT ATA AAG TTA GTA ATT TTA ACT ATA TCA TTT GAT ACA TTA TTA     274
 H   I   F   I   K   L   V   I   L   T   I   S   F   D   T   L   L

GGA TGT TTA AGT GCA ATA AAA AGT CGT AAA TTT AAT TCT AGT TTT GGA ATA     325
 G   C   L   S   A   I   K   S   R   K   F   N   S   S   F   G   I

GAT GGA GGA ATC AGA AAA GTA GCA ATG ATA GCA TGT ATA TTT TTT TTA TCA     376
 D   G   G   I   R   K   V   A   M   I   A   C   I   F   F   L   S

GTA GTT GAC ATT CTT ACA AAG TTT AAC TTT TTA TTT ATG TTA CCA CAA GAT     427
 V   V   D   I   L   T   K   F   N   F   L   F   M   L   P   Q   D

TGT ATC AAT TTT TTA AGA CTA AAA CAT CTT GGA ATA TCT GAA TTT TTC TCT     478
 C   I   N   F   L   R   L   K   H   L   G   I   S   E   F   F   S

ATT TTA TTT ATT TTA TAT GAA AGT GTA AGT ATA TTA AAA AAT ATG TGC TTA     529
 I   L   F   I   L   Y   E   S   V   S   I   L   K   N   M   C   L

TGT GGA TTA CCA GTA CCT AAG AGA TTA AAG GAA AAA ATA GCA ATT TTA CTA     580
 C   G   L   P   V   P   K   R   L   K   E   K   I   A   I   L   L

GAT GCA ATG ACA GAT GAA ATG AAT GCT AAG GAT GAA AAG TAA GTAATGGT        630
 D   A   M   T   D   E   M   N   A   K   D   E   K  END

AGATATAATA AAGATATTAA CAAATAAAAA GTGTTATCCA AATAAGAATA GCTGAAAGTT       690

ATCATAATTC ATGAAACTAA TAATGAAAAC GAGGGAGCAG ATGCCAAGAG ACACACAAGT       750

ATTAAATACA TATAATTTCG AAGCAAGTGT TCATTACTAT ATAGATGACA AGGTAGTATA       810

TCAAACATTG GTTCACAAAG ATGGTGCATG GTCAGTTGGT AAAATCTATT AAGCTACATT       870

AGTTACAGAT ATCACAAACT ATAATAGTTA AACATAGAAA TATGTGTAAA TTGTGATGGA       930

AATTATTCAA AAACACAAAA ATACGTGATG AAGGACAAAA TGATATAGAA AATAAGTATC       990

AAACCTTAAT AAATGATTTA ATTGATAGTT TAAAAGTTAT AGGAAAAATA TATAAAGAAA      1050

TAAAAACATT AAAAAAAATAT AAGATATGTT TACAAATTAC TATCAGACAA TCTCCTTATC     1110

TAATAGAAGA GTCAATTAAC TAATTGAGTA TCTTTAAATT GAAATGTTAG GAAGTGATTT      1170

AAATATGAAA ACTTAAATT  1189
```

FIG. 8. Nucleotide and deduced amino acid sequences of the small open reading frame located between the 3' end of *C. difficile* toxin B and the 5' end of toxin A. Also included are the sequences between toxin B and the open reading frame and between the open reading frame and the first nucleotide (−160) listed in the toxin A sequence (Fig. 1).

surface proteins of *Plasmodium* species have repeated sequences, several of which are believed to be target cell binding proteins (25). These repeating units range from 3 to 18 amino acids in length, are repeated from 5 to as many as 41 times, and may consist of nearly 40% of the protein (5, 6). Several toxin genes have been sequenced that contain repeating sequences at the C-terminal end of the proteins. The C-terminal region of the *E. coli* hemolysin polypeptide contains 13 8-amino-acid repeating units, which are required for hemolytic activity (9). The calmodulin-sensitive adenylate cyclase of *Bordetella pertussis* contains two regions that contain repeating units (10). Eleven repeating units of 15 amino acids have recently been reported for the insecticidal crystal proteins of *Bacillus thuringiensis* (34). Although the repeating sequences of *C. difficile* toxin A do not have any sequence similarities with any of these other proteins, location at the C-terminal end of the proteins is common, and some may have a common role for target cell binding. Because the repeating region constitutes about one-third of the entire toxin molecule and the repeats are highly hydrophilic, it would be interesting to determine the spatial distribution of these repeats in the native protein. It remains to be shown whether a periodicity on the surface of the toxin molecule confers certain unique biological properties to the protein. We are currently pursuing research in this area to gain more understanding of the structure and function of this toxin.

## LITERATURE CITED

1. Banno, Y., T. Kobayashi, K. Watanabe, K. Ueno, and Y. Nozawa. 1981. Two toxins (D-1 and D-2) of *Clostridium difficile* causing antibiotic-associated colitis: purification and some characterization. Biochem. Int. 2:629–635.

2. Bartlett, J. G., N. S. Taylor, T. Chang, and J. Dzink. 1980. Clinical and laboratory observations in *Clostridium difficile* colitis. Am. J. Clin. Nutr. 33:2521–2526.

3. Biggin, M. D., T. J. Gibson, and G. F. Hong. 1983. Buffer gradient gels and 35S label as an aid to rapid DNA sequence determination. Proc. Natl. Acad. Sci. USA 80:3963–3965.

4. Birnboim, H. C., and J. Doly. 1979. A rapid alkaline extraction procedure for screening recombinant plasmid DNA. Nucleic Acids Res. 7:1513–1523.

5. Coppel, R. L., A. F. Cowman, K. R. Lingelbach, G. V. Brown, R. B. Saint, D. J. Kemp, and R. F. Anders. 1983. Isolate-specific S-antigen of *Plasmodium falciparum* contains a repeated sequence of eleven amino acids. Nature (London) 306:751–756.

6. Dame, J. B., J. L. Williams, T. F. McCutchan, J. L. Weber, R. A. Wirtz, W. T. Hockmeyer, W. L. Maloy, J. D. Haynes, I. Schneider, D. Roberts, G. S. Sanders, E. P. Reddy, C. L. Diggs, and L. H. Miller. 1984. Structure of the gene encoding the immunodominant surface antigen of the sporozoite of the human malaria parasite *Plasmodium falciparum*. Science 225:593–599.

7. Ehrich, M., R. L. Van Tassell, J. M. Libby, and T. D. Wilkins. 1980. Production of *Clostridium difficile* antitoxin. Infect. Immun. 28:1041–1043.

8. Eichel-Streiber, C. V., D. Suckau, M. Wachter, and U. Hadding. 1989. Cloning and characterization of overlapping DNA fragments of the toxin A gene of *Clostridium difficile*. J. Gen. Microbiol. 135:55–64.

9. Felmlee, T., and R. A. Welch. 1988. Alterations of amino acid repeats in the *Escherichia coli* hemolysin affect cytolytic activity and secretion. Proc. Natl. Acad. Sci. USA 85:5269–5273.

10. Glaser, P., D. Ladant, O. Sezer, F. Pichot, A. Ullmann, and A. Danchin. 1988. The calmodulin-sensitive adenylate cyclase of *Bordetella pertussis*: cloning and expression in *Escherichia coli*. Mol. Microbiol. 2:19–30.

11. Hanahan, D. 1983. Studies on transformation of *Escherichia coli* with plasmids. J. Mol. Biol. 166:557–580.

12. Holdeman, L. V., E. P. Cato, and W. E. C. Moore (ed.). 1977. Anaerobe laboratory manual, 4th ed. Virginia Polytechnic Institute and State University, Blacksburg.

13. Johnson, J. L. 1981. Genetic characterization, p. 450–472. *In* P. Gerhardt, R. G. E. Murray, R. N. Costilow, E. W. Nester, W. A. Wood, N. R. Krieg, and G. B. Phillips (ed.), Manual of methods for general bacteriology. American Society for Microbiology, Washington, D.C.

14. Johnson, J. L., C. H. Dove, S. B. Price, T. W. Sickles, C. J. Phelps, and T. D. Wilkins. 1988. The Toxin A gene of *Clostridium difficile*, p. 115–123. *In* J. M. Hardie and S. P. Borriello (ed.), Anaerobes today. John Wiley & Sons, Ltd., London.

15. Krivan, H. C., G. F. Clark, D. F. Smith, and T. D. Wilkins. 1986. Cell surface binding site for *Clostridium difficile* enterotoxin: evidence for a glycoconjugate containing the sequence Galα1-3Galβ1-4GlcNAc. Infect. Immun. 53:573–581.

16. Lyerly, D. M., H. C. Krivan, and T. D. Wilkins. 1988. *Clostridium difficile*: its disease and toxins. Clin. Microbiol. Rev. 1:1–18.

17. Lyerly, D. M., C. J. Phelps, and T. D. Wilkins. 1985. Monoclonal and specific polyclonal antibodies for immunoassay of *Clostridium difficile* toxin A. J. Clin. Microbiol. 21:12–14.

18. Lyerly, D. M., M. D. Roberts, C. J. Phelps, and T. D. Wilkins. 1986. Purification and properties of toxins A and B of *Clostridium difficile*. FEMS Microbiol. Lett. 33:31–35.

19. Lyerly, D. M., K. S. Saum, D. K. McDonald, and T. D. Wilkins. 1985. Effects of *Clostridium difficile* given intragastrically to

animals. Infect. Immun. **47**:349–352.

20. **Lyerly, D. M., N. M. Sullivan, and T. D. Wilkins.** 1983. Enzyme-linked immunosorbent assay for *Clostridium difficile* toxin A. J. Clin. Microbiol. **17**:72–78.

21. **Maniatis, T., E. F. Fritsch, and J. Sambrook.** 1982. Molecular cloning: a laboratory manual. Cold Spring Harbor Laboratory, Cold Spring Harbor, N.Y.

22. **Marmur, J.** 1961. A procedure for the isolation of deoxyribonucleic acid from microorganisms. J. Mol. Biol. **3**:208–218.

23. **Matsudaira, P.** 1987. Sequence from picomole quantities of proteins electroblotted onto polyvinylidene difluoride membranes. J. Biol. Chem. **262**:10035–10038.

24. **Michaelis, S., and J. Beckwith.** 1982. Mechanism of incorporation of cell envelope proteins in *Escherichia coli*. Annu. Rev. Microbiol. **36**:435–465.

25. **Miller, L. H., R. J. Howard, R. Carter, M. F. Good, V. Nussenzweig, and R. S. Nussenzweig.** 1986. Research toward malaria vaccines. Science **234**:1349–1356.

26. **Muldrow, L. L., G. C. Ibeanu, N. I. Lee, N. K. Bose, and J. Johnson.** 1987. Molecular cloning of *Clostridium difficile* toxin A gene fragment in lambda gt11. FEBS Lett. **213**:249–253.

27. **Pothoulakis, D., L. M. Barone, R. Ely, B. Faris, M. E. Clark, C. Franzblau, and J. T. LaMont.** 1986. Purification and properties of *Clostridium difficile* cytotoxin B. J. Biol. Chem. **261**:1316–1321.

28. **Price, S. B., C. J. Phelps, T. D. Wilkins, and J. L. Johnson.** 1987. Cloning of the carbohydrate-binding portion of the toxin A gene of *Clostridium difficile*. Curr. Microbiol. **16**:55–60.

29. **Rihn, B., J. M. Scheftel, R. Girardot, and H. Monteil.** 1984. A new purification procedure for *Clostridium difficile* enterotoxin. Biochem. Biophys. Res. Commun. **124**:690–695.

30. **Rolfe, R. D., and S. M. Finegold.** 1979. Purification and characterization of *Clostridium difficile* toxin. Infect. Immun. **25**:191–201.

31. **Sanger, F., S. Nicklen, and A. R. Coulson.** 1977. DNA sequencing with chain-terminating inhibitors. Proc. Natl. Acad. Sci. USA **74**:5463–5467.

32. **Sullivan, N. M., S. Pellett, and T. D. Wilkins.** 1982. Purification and characterization of toxins A and B of *Clostridium difficile*. Infect. Immun. **35**:1032–1040.

33. **Taylor, N. S., G. M. Thorne, and J. G. Bartlett.** 1981. Comparison of two toxins produced by *Clostridium difficile*. Infect. Immun. **34**:1036–1043.

34. **Widner, W. R., and H. R. Whiteley.** 1989. Two highly related insecticidal crystal proteins of *Bacillus thuringiensis* subsp. *kurstaki* possess different host range specificities. J. Bacteriol. **171**:965–974.

35. **Wilkins, T. D.** 1987. Role of *Clostridium difficile* toxins in disease. Gastroenterology **93**:389–391.

36. **Wren, B. W., C. L. Clayton, P. P. Mullany, and S. Tabaqchali.** 1987. Molecular cloning and expression of *Clostridium difficile* toxin A in *Escherichia coli* K12. FEBS Lett. **225**:82–86.

37. **Yanisch-Perron, C., J. Vieira, and J. Messing.** 1985. Improved M13 phage cloning vectors and host strains: nucleotide sequences of the M13mp18 and pUC19 vectors. Gene **34**:103–119.