

Cloning and chromosomal location of human $\alpha 1(\text{XVI})$ collagen

TE-CHENG PAN*, RUI-ZHU ZHANG*, MARIE-GENEVIEVE MATTEI†, RUPERT TIMPL‡, AND MON-LI CHU*§

*Departments of Biochemistry and Molecular Biology, and of Dermatology, Jefferson Institute of Molecular Medicine, Thomas Jefferson University, Philadelphia, PA 19107; †Center de Genetique Medicale, Institut National de la Santé et de la Recherche Médicale U242, Marseille, France; and ‡Max-Planck-Institut für Biochemie, D-8033 Martinsried, Federal Republic of Germany

Communicated by Darwin J. Prockop, March 24, 1992

ABSTRACT We have characterized cDNA clones that encode a newly discovered collagenous polypeptide. A 4-kilobase (kb) cDNA clone was initially isolated by screening a human fibroblast cDNA library with a probe encoding the collagenous domain of the human $\alpha 3(\text{VI})$ collagen. Subsequent screening of another fibroblast cDNA library yielded overlapping clones having a total length of 5.4 kb, which contained an open reading frame of 1603 amino acids including a 21-amino acid signal peptide. The predicted polypeptide consists of 10 collagenous domains 15–422 amino acids long, which were interspersed with 11 noncollagenous (NC) domains. Except for a large N-terminal NC11 domain of 312 residues, most of the NC domains were short (11–39 residues) and cysteine-rich. The overall structural arrangement differed significantly from other known collagen chains. Further analysis indicated that the deduced polypeptide exhibited several structural features characteristically seen in members of the fibril-associated collagen, types IX, XII, and XIV. In addition, the cysteine-rich motifs in the NC domains resembled those found in the cuticle collagen of *Caenorhabditis elegans*. Northern blot analysis showed hybridization of the cDNA to a 5.5-kb mRNA in human fibroblasts and keratinocytes. The gene was localized by *in situ* hybridization to band p34–35 of human chromosome 1. The data clearly support the conclusion that the cDNA encodes a collagen chain that has not been previously described. We suggest that the cDNA clones encode the $\alpha 1$ chain of type XVI collagen.

Collagens, the major constituents of connective tissues, represent a large family of structurally related proteins with distinct tissue distributions and functions (see refs. 1 and 2 for recent reviews). To date, 14 different types of collagen have been described in vertebrates and have been assigned Roman numerals. Based on their primary structure and supramolecular assembly, collagens can be divided into two major classes: the fibril-forming collagens and the non-fibril-forming collagens. A long central triple-helical domain, without Gly-Xaa-Xaa interruptions, is the hallmark of the former class. Types I, II, III, V, and XI, which form highly organized fibrils in a quarter-staggered fashion, are members of this class. The remaining types belong to the latter class. These collagens are very heterogeneous in size, but a common feature is the presence of imperfections in the Gly-Xaa-Xaa repeating pattern. Within the latter class, types IX, XII, and XIV collagens form a subgroup named the fibril-associated collagens with interrupted triple helices (FACIT) (3). These collagens are associated with type I or II collagen fibrils, which play a role in interaction of these fibrils with other matrix components (4–6). Members of the FACIT group share several common structural features, although their sizes and primary structures vary greatly. Some of these characteristic features are also found in the cuticle collagens of *Caenorhabditis elegans* (7).

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.

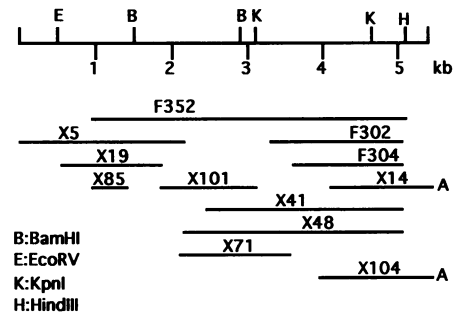


FIG. 1. Partial restriction map of the composite human $\alpha 1(\text{XVI})$ collagen cDNAs. Restriction enzymes used were *Bam*HI (B), *Eco*RI (E), *Kpn*I (K), and *Hind*III (H). Clones were isolated from cDNA libraries prepared from either fibroblast 3349 (F) or fibroblast 1262 (X). A, poly(A) tail.

In this report, we describe the isolation and characterization of cDNA clones that encode a newly discovered collagenous polypeptide. We suggest that the putative collagenous polypeptide be named the $\alpha 1$ chain of type XVI collagen.[¶]

MATERIALS AND METHODS

Cell Strains. Human skin fibroblasts 3348 and 3349 were obtained from the Human Genetic Mutant Cell Repository, Camden, NJ. Human fetal skin fibroblasts 1106 and 1262 were obtained from the American Type Culture Collection. Human epidermal keratinocytes were purchased from Clonetics, San Diego.

Isolation and Sequencing of cDNA Clones. Two human fibroblast cDNA libraries were used to screen for overlapping clones. One library was constructed by using the poly(A)⁺ RNA from fibroblast 3349 and the oligo(dT) primer as described (8). The other library was prepared from poly(A)⁺ RNA isolated from fibroblast 1262, a patient with a perinatal lethal osteogenesis imperfecta (9). The library was prepared from 5 μg of poly(A)⁺ RNA using a cDNA synthesis kit with mixed hexanucleotides as primers (Pharmacia), and the resultant cDNAs were ligated and packaged into the vector λ ZAP II (Stratagene). The unamplified cDNA libraries were screened with cDNA fragments labeled with ³²P by nick-translation (10) in a hybridization mixture containing 5 \times SSC (1 \times SSC is 0.15 M NaCl/0.015 M sodium citrate, pH 7.0) at 60°C. The filters were washed at 60°C in 0.5 \times SSC.

cDNA inserts were sequenced by the dideoxynucleotide chain-termination method (11) using adenosine 5'-[γ -³⁵S]thio]triphosphate (New England Nuclear) and the mod-

Abbreviations: COL domain, collagenous domain; NC domain, noncollagenous domain; FACIT, fibril-associated collagens with interrupted triple helices.

[§]To whom reprint requests should be addressed at: Department of Biochemistry and Molecular Biology, Thomas Jefferson University, 233 South 10th Street, Philadelphia, PA 19107.

[¶]The sequence reported in this paper has been deposited in the GenBank data base (accession no. M92642).

GGGACCCGTGACCGTACGGAGAGATAGAGAACCCCTAGCGGACCTGAGGAACCTGGGCT 60
 CTGCTTCCCTCCGCTGCTGCTGAGTCTCACTGAGTGGGGGATCCCGGGCCACCG 120
 TGCTGGCTGTAGCTGACCTTTTGACCCGGATGGGTATCTGGGCTCTGGGCTG 180
 M W V S W A P G L 9
 TGCTGCTGGTCTTGGGCTACCTTGGGCAATGGGCAATACAGTGCACAAATGCCA 240
 L L G L L M W A T A T G G C C A G A N T G A Q (C) P 29
 CCTTACACAGGAGGACTCAAAATGGAACACAGTACTAGCTCCAGCCAACTGACT 300
 P S Q Q E G L K L E H S S S L P A N V T 49
 GGCTCAAOCTTCCAGCGACTCAGCTCAGTGAAGTCTGCATCAAGAAGATCCG 360
 G F N L I H R L S L M K K S A I K K I R 69
 AACCCAAAGGGCTCTCATCTCGGCTGGGGGCGCCCGTACCAGCCAGCCACCGA 420
 N P K G P L L R L L G A A P V T Q P T R 89
 CGAATTTCCCTGGGCTCTCGGAGAACTTCCCTGTGCTGACTACTLCTCAAGR 480
 R V F P R G L P E E F A L V L L T L L L K 109
 AAACACACCCAGAGAGCTGTATCTGTTCAAGTGACCGATGCAATGGGTATCCA 540
 K H T H Q K T W Y L F Q V T D A N G Y P 129
 CAGATTCCTGGAGTCAACAGCAAGCCAGCTGAGCTGAGCCAGGGCCAGGCGAG 600
 Q I S L E V N S Q E R S L E L R A Q G Q 149
 GATGGCACTTTGCTGCTGACTTCCAGTCCCGAGCTCTTGCATTCGCTGGGAC 660
 D G D F V S (C) I F P V P Q L F D L R W H 169
 AACCTGATCTGGAGTGGCTGAGCTGGCTGTGTCAGCTGAGTCTGAGCTCAGCC 720
 L M L S V A G R V A S V H V D (C) S S A 189
 TCCTCCAGCTCTGGGCGCCAGCACCATGAGGCTGTGGGCATGTATTCTAGCC 780
 S S Q P L G P R P M R P V G H V F L G 209
 TTGAGCTGAGCAGCGAACTGCTGCTTGCCTTACCTTACAGGCTGACATCTACTGT 840
 L D A E Q G K P V S F D L Q Q V H I Y (C) 229
 GACCCGAGCTGCTGCTGAGAGGGCTGCTGTGAGATTTACAGCAGGCTCCGCCCA 900
 D P E L V L E E G (C) I L P A G (C) P P 249
 GACAGCTCAAGCCCGCCGACCCAGCAAGTGAAGCTATTAGATCAATCCAGAG 960
 E T S K A R R R D T Q S N E L I E I N P Q 269
 TCTGAAGCAAGCTTACACCCGCTCTTCTGCTGGAGGAGCCAAAACAGCGAGT 1020
 S E G K V Y T R (C) L E E P Q N S E 289
 GATGCCAGCTGACCGAAGATCAGCCAGAGCAGAAAGGGAGCAAGGTCCATCAG 1080
 D A Q L T G R I S Q K A E R G A K V H Q 309
 GAGACAGCCGATGAGTCTCCGCTGCTGCTGAGTGGCCGGGACAGCAATGTACA 1140
 E T A A D E (C) P P P V H G A R D S N V T 329
 CTGCTCCCTTGGCCCAAGGAGGAAAGTGGGGGCTGCTGCTGCTCACCAGGG 1200
 L A P S P K P G K G E R G L P G P P G 349
 TCCAAGGGAGAGAGGAGCCAGGGCAATGAGTGTGGAATCTCCCGGATGCCCG 1260
 S K G E K G A R G N D (C) V R I S P D A P 369
 CTTAGTGTGAGAGAGGAGGAGGAGGAGCTCAGGAGCTTGGAGCTCA 1320
 L Q (C) A E G P K G E K G E S G A L G L P S 389
 GGACTCCAGGCTCAACAGGCGAAGGGCCAGAAAGGGAGAAAGGGAGGAGGATC 1380
 G L P G S T G E K G K G K G K G D G G I 409
 AAGGGCTCGGGAAAGCCAGGGCCAGCCAGGAGAGATCTGTCTATTGGGCC 1440
 K G V P G K P G R D A P E I (C) V I G P 429
 AAAGGCGAAGGAGACCTTGGCTTGTGGGCTGAGGGCTGAGGAGAGCTGGG 1500
 K G Q K G D P G F V G P E G L A G E P G 449
 CCCCAGGCTCCCTGAGCCCTGGATAGGACTGCTGGACCCGGGGATCCAGGT 1560
 P P G L P G P P P (C) I G L P G T P G D P G 469
 GGCCACAGGCCCCAAGGAGCAAGGGCAGCTGGGATCCCGAAAGGAAGGCC 1620
 P P P G P K G D K G S G I P G K E G P 489
 GGTGGAACTGGGAAGCCAGGTGTGAAGGAGAGAGGGTGAOCCCTGTAAGTGTG 1680
 G K P G K P G V K G E K G D E (C) E V (C) 509
 CCAACACTGCTGAAGGTTCCAGAATTTGTGACTTCTGGAAGCCAGGGCCAAA 1740
 P T L P E G F Q N F I G L P G K P G P K 529
 GGGAGCTGTGATCTGTGAGGAGGAGGAGCCCTGGCATCAAGGCATCAAGGA 1800
 G E P G D P V R A R G G D P G I Q G I K G 549
 GAGAAGGGAGGAGGCTGCTGCTGCTGCTGAGTGTAGGGCCAGCATCTTGTG 1860
 E K G E P (C) L S S V V G A Q H L V S 569
 TCAAGGGCCAGTGGAGATGGGTTCCCTGCTTGTGCTGCTGGCTCCCGGT 1920
 S T I G A S G D V G S P P E G L P G L P G 589
 AGAGCTGGGTTCCAGGCTGAAGGAGAGAGGTAATCTCGGGAGGAGGCGGAGCT 1980
 R A G V P G L K G E K G N P G E A G P A 609
 CGCAGCTCCAGGGCCAGGAGCAGTGGGCGCCAGGACTCAAGGGCCAGGGGAG 2040
 G S P P P P P G P V G P A G I K G A K G E 629
 CDTGTGAGCCCTGCGAGCTGTCACCTTCAAGTGGGATGCTCGGTGGTGGCC 2100
 P (C) E P (C) P A L S N L Q D G D V R V A 649
 TTGCTGGCCATCCGAGAGAGGGGAACTGGGCTTCCAGCTTGGCTTGGCCAGCA 2160
 L P G P S G E K G E P P G E G L P G 669
 AAACAGGCGAGCTGAGAGCTGAGCTGAAGGGCAGAGGGTGAATGCTGGAACTCT 2220
 K Q G K A G E R G L K G Q K G D A G N P 689
 GGAGACCTTGAACCGGGCCAGCCAGGCGCCAGGACTGACAGAGGCTGGAGTT 2280
 G D P P G T P G T I T A R R P G L S G E P G V 709
 CAGGCGCCGGGGCCAAAGGAGAAAGGTTATGGTCACTTCTCCCGAGCTG 2340
 Q G G P A G P K G E K G D G (C) T A (C) P S L 729
 CAGGGAGCAGTACAGACATGGCAGGAGCCCTGGGACCCCGGCAAGGAGAACAG 2400
 Q G T V T D M A G R P P G Q P P K G E Q 749
 GGGCCAGGAGGCTGGCCAGCTGTAACCGGGCAACCGGCTTACCAGGAGTTCAG 2460
 G P E G V G R P G K P G Q G L P G V Q 769
 GGGCCCCAGGACTGAAGGGTGCAGGAGAGCCAGGCTCCAGGAAGGAGTCCAG 2520
 G P P G L K G V Q G E P G P P E R G V Q 789
 GGACCCAGGGAGCTGGAGCCCGGGTTTCCGCTTCCAGGACTTCCGGGACT 2580
 C P Q G E P P G A P G L P G I Q G L P G P 809
 CGGGACACCTGGCCACTGAGAGAGGGTCCCGAGGATCTCCAGGGTGAAGGA 2640
 R G P P G P T G E K G A Q G S P G V K G 829
 GCCACGGACCCCTGGGACTCTGGCCAGTCTCTGGGCTCCGGCCGTGATGG 2700
 A T G P V G P P G A S V S G P P G R D G 849

COL 10

COL 9

COL 8

COL 7

COL 6

COL 5

CAGCAAGGACAGCGGACTCAGAGGAACACAGGTGAAAGAGACCGAGGAGAGAG 2760 COL 5
 Q Q G Q T G L R G T P G E K G P R G E K 869
 GGTGAGCCAGGGAGTCTCTCCCTCCCTCAAGGAGACCTTCTCTTGGCATGCC 2820
 G E P G E (C) S (C) P S Q G D L I F S G M P 889
 GGTGCTCGGACTTTCAGTGGCAGCTTGGCAGCCGGGCGCGGCTCCACAGGT 2880
 G A P L W M G S S W O F G P P G 909 COL 4
 ATTCCGAGCAGCAGCGCCCTCCGAGTACCTGGCTGCGAGGAGTCCCTGGAACAC 2940
 I P P G P G P P G V P P G L Q G V P G G N N 929
 GGTTTGCCAGGACGCTGGGCTCAGCAGAACGGGATCTTACCAATTGAACAGCAC 3000
 G L P G Q P G L T A E L G S L F I E Q H 949
 CTCCTAAGAGTATCTCGGGGACTGTGTCAGGGCAGAGGGCCACCCAGGTACCT 3060
 L L K S I (C) G D (C) V Q G Q R A H P G Y L 969
 GTGGAGAGGAGAGGAGAGCAGGCACTCCCTGGTGTCCAGGCTCCGAACTG 3120
 V E K G E K G D Q G I P G G L G D N (C) 989 COL 3
 GCCACTGCTTTTGTACTGGAGCCCAAGAGCCAGGGGCGGGTGCACACAGT 3180
 A Q (C) F L S L E R P P R A E E A R G D N S 1009
 GAGGAGATCTGGCTGTTGGAGCCAGGCTTCCCTGCTCCGGATTCGCCAGCC 3240
 E G D P G (C) V G S P G L P G P P G L P G 1029
 CAGAGAGGAGAGGGTCCGCTGGCAGGAGGCTCCCGGCTCCAGCCCTATC 3300
 Q R G E E G P P G M R G S P G P P P G P I 1049
 GGCCCCAGGTTTCTGGTGTGTGCTCCCGGATTCGCTGGCTCAAGGAGAG 3360
 G P P G F P G A V G S P G L P G G E 1069
 CGAGCTCAGGGCTGACTGAGCAAGGGGGAGCGGCTCCAGGGCAACAGGT 3420
 R G L T G L T G D K G E P G Q P G G 1089
 TACCGAGTGCAGGCGCCCGGAGTCCCTGCATCAAGGGGAGGCTGTGTACAC 3480
 Y P G A T G P P G L P G I K G E R G Y T 1109
 GGTGAGGGAGAGAAAGGAGAGCCGGCCAGGACTTGAAGCCCTCCAGGCC 3540
 G S A G E K G E P G P P G S E G L P G P 1129
 CCAGCCAGGGTCCAGAGGAGGAGCCAGGACCCAGTACTCCGTGAGAGGGC 3600
 G P P A G C P R G E R G P Q G N S G E K G 1149
 GACCAGGATTCAAGCCAGGCTTACCGGSCACCGGTCCCTGGATCCCA 3660
 D Q G F Q G Q P G F T G P T G P T 1169
 GGCAAGTTGATCAGCTGGCCAGCTGCGCTCAGGAGAAAGCCAGGAGGATT 3720
 G K V G S P G P P G P G P (C) A E K G S E G I 1189
 CAGGCCATCAGGCTGCTGCTCCCTGGCCAGCCGACTCCCTGGGATCAGGCC 3780 COL 2
 R G P S G L P G S P G P P G P G I Q G 1209
 CCGCGGCTTGTGATGTTGGATGGAGAGCCAGGCTGGCTTCCAGGGAGCCT 3840
 P A G L D G L D G K D G K P G L R G D P 1229
 GGTCTGTGCCCCCTGGACTTGGACCCAGGCTTTAAGGGAAGAAACAGGACAT 3900
 T A G P P G L M G P P G P K G T G H 1249
 CTTGCTCCAGGACTAAGGCTGCTGGAAGCCAGGCTCCCTGCGAGCATGCC 3960
 P L G P G P G D (C) G K P G P P G S T G 1269
 CGGCTGGCCAGAGGCTGAACCTGGTGGTCCAGCCAGGAAAGCCCGTCCCG 4020
 R P G A E G E P P A M G P Q G R P P G P 1289
 GGACAGTTGGCCAGCAGGCTCCAGCCAGGAGGAGCCAGCCGATGCTGTGCA 4080
 G H V G P P G P P G Q P G P A G I S A V 1309
 GGTCTGAAAGAGCCAGGAGCCAGGAGAAAGGGCCCTGAGCCCTCCAGCC 4140
 L K G D R G A T G E R G L A G L P G Q 1329
 CCGGCCCCGAGGAGCTGCGCCAGGAGCTGGTACGATGCTGAGCTGGC 4200
 P P P G H P P P G P G E P A G T D G A A G 1349
 AAAGAGGACCCCTGGAAGCAGGATTTATGACCTCTGCTCCAGGATCCA 4260
 K E G P P G K Q G F Y G P P G P K G D P 1369
 GGAGCTCAGCAGAGGAGCCAGCAGAGAGGAGCCAGCCGATGCTGTGCA 4320
 G A A G Q K G Q A G E K G R A G M P P G 1389
 CCTGCAAGATGGTTCATGGGCTGTGGCCAGCCGGCCCTGAGGAGAGAGCC 4380
 P G K S G S M G P V G P P G P A G E R G 1409
 CAGCTGAGCTCCGGCCCTGGGAGCCCTGGTCTGCTGCTCCCTGCTGCT 4440
 H P G A P G P S G S P G L P G V G S M 1429
 GGAGACATGTAATTATGATGAATCAAGAGTTTATCAGCAAGAGATCAATAA 4500
 G D M V N Y D E I K R F I R Q E I I K M 1449
 TTTGATGAGAGATGGCTTACTACCTCCAGGATGCAATCCCAATGAGATGGCG 4560
 F D E R M A Y Y T S R M Q F P M E M A A 1469
 GCTCCGAGCAGGCGCTCAGGAGAGATGGTCTCCGGCAGCCAGTGTCCA 4620
 A P G R P P P P G K D G A G P R P G A P 1489
 GGTGCTGGCTCCCTGGTCAATGCGAGAGGAGGAGGAGGCTTGGCAGGAGTA 4680
 G S P G L P G Q I G R E G R Q G L P G V 1509
 AGAGGATTCCTGGTCAAAAGGTGAAAAGGGGACATGGTATTGGATTCAGGAGAA 4740 COL 1
 R G L P G T K G E K G D I G I G I A G E 1529
 AATGGCTTCCGGCCCCAGGCTCCTCAAGTCTCCAGCTATGCGAGGATGGTGA 4800
 N G L P G P P P G P Q G P P G Y G K M G A 1549
 ACAGGACCAATGGCCAGCAAGCATCCCTGGCATCCCTGGCCCGGGTCCCATGG 4860
 T G P M G Q Q G I P G I P G P P G P M G 1569
 CAGCAGGCAAGGCTGCGACTTAACTCCCTGCTTGGGCTTCCAGGATGGAG 4920
 Q P G K A G H (C) N P S D (C) F G A M P M E 1589
 CAGCAGTACCCACCAAGAAAATGAAAGGGGCTTTGGTGAATTCACCACTGCC 4980
 Q Q Y P P M K T M K G P F G 1603
 TTTGATGAAAGACTCGTTGGGAATAAATGGCCAAAGCTTATAGGACTCTGTGACAG 5040
 TGTGAATGTTTTTTTTTTGTTGTTGTTTTTAAATGCTGTTAATTTTTAAATA 5100
 A T A A A A A A A A A A T A T C C C T T C C T T C C A G T G G T T C C T G T G C T G A C C A 5160
 G A G C T C C T T G C C T C T T T C C G T T A G T C C A G A A A A A G G C A T T T G G T A 5220
 C C C G C A T A C T G T A T C T A G C A T T C A A G G G C T G A G T G G G A G C T G C T G 5280
 C A G G G A T T T A G A C C A T P A G C A C A T A A T A G C A G T C T C G A A A A A A A A 5340
 T T A A A A A A A A G G G T A T T C C C T C T T A A A A A A A A A A A A A A 5387

Fig. 2. Nucleotide and deduced amino acid sequence of human $\alpha 1(XVI)$ cDNA. Boxes, collagenous (COL) domains. Cysteines are circled. Imperfections in the Gly-Xaa-Xaa repetitive pattern, polyadenylation signals (AATAAA), N-glycosylation sites (Asn-Xaa-Thr/Ser), and Arg-Gly-Asp tripeptides are underlined. Arrow, potential cleavage site for signal peptidase.

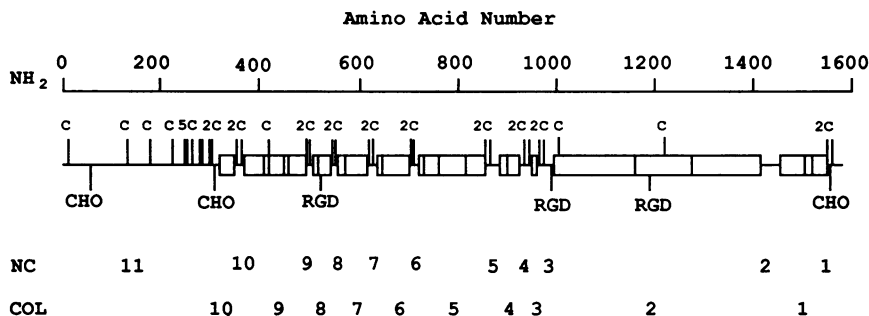


FIG. 3. Schematic of the domain structure of human $\alpha 1(\text{XVI})$ collagen, which consists of 10 COL domains (open boxes) and 11 NC domains (horizontal lines). Imperfections of Gly-Xaa-Xaa repeats are indicated by vertical lines within open boxes. Positions of cysteines (C), N-glycosylation sites (CHO), and Arg-Gly-Asp (RGD) tripeptides are indicated.

ified T7 polymerase (Sequenase kit, United States Biochemical). The sequences were analyzed by the PC/GENE computer program (IntelliGenetics). GenBank searches and the alignment of amino acid sequences were performed by using the FASTA program (12).

Isolation of RNA and Northern Blot Analysis. Total RNAs were isolated from cultured cells by acid guanidine thiocyanate/phenol/chloroform extraction (13). Poly(A)⁺ RNAs were selected on an oligo(dT)-cellulose column (Collaborative Research). RNA samples were electrophoresed on a 1% agarose gel containing 6% formaldehyde, transferred to nitrocellulose (14), and hybridized to ³²P-labeled cDNAs as described (15). The cDNA probes for the $\alpha 1(\text{I})$, $\alpha 1(\text{III})$, and $\alpha 1(\text{VI})$ collagen chains were Hf677 (16), Hf934 (17), and P18 (18), respectively.

Chromosomal Localization. Human metaphase chromosomes prepared from phytohemagglutinin-stimulated lymphocytes of a normal male individual were used for chromosomal *in situ* hybridization as described (19). The cDNA probe used, F352, was labeled with [³H]dCTP and [³H]dTTP by nick-translation (10).

RESULTS

Isolation of cDNA Clones. A 4-kb cDNA clone, F352, was isolated from the human fibroblast 3349 cDNA library by screening 5 × 10⁵ clones with a 1.5-kilobase (kb) cDNA, P24, which encodes the $\alpha 3$ chain of human type VI collagen (20). Partial DNA sequence analysis of the cDNA indicated that it encoded a collagenous polypeptide that had not been previously described. Rescreening of the cDNA library with F352 yielded only two shorter clones, F302 and F304 (Fig. 1). To obtain cDNA clones corresponding to the full-length mRNA, a random-primed cDNA library was constructed by using poly(A)⁺ RNAs from fibroblast 1262. Screening of 5 × 10⁵ independent plaques yielded 9 clones. DNA sequencing revealed that these clones overlapped and together covered ≈5.4 kb of the mRNA.

Nucleotide and Amino Acid Sequences. Nucleotide sequence analysis indicated that one of the clones (X5) began in the 5' untranslated region of the mRNA, which was followed by the sequence encoding a putative 21-amino acid signal peptide. Two of the clones, X14 and X104, ended in a poly(A) tail. Therefore, these overlapping clones cover essentially the full-length mRNA. The cDNA clones encode an open reading frame of 1603 amino acids, which starts with an ATG codon at nucleotide 154 and ends with a TGA stop codon at nucleotide 4963 (Fig. 2). Following the stop codon

there are 425 base pairs of the 3' noncoding sequence, which contains 4 consensus polyadenylation signals (AATAAA).

Analysis of the predicted protein product revealed the presence of 10 collagenous (COL) domains separated by 11 noncollagenous (NC) regions (Figs. 2 and 3). The COL domains range in size from 15 to 422 amino acids; most of them contain Gly-Xaa-Xaa imperfections (Table 1). The NC domains are all relatively short (11–39 amino acids), except for the large N-terminal NC11 domain of 312 amino acids. The most notable feature is the presence of numerous cysteines (a total of 32) in the molecule, almost all of which are found in the NC domains. A majority of the NC domains contain two cysteines separated by two other amino acids, which are often located at the end of the previous COL domains. Of particular interest is the sequence at the end of the COL1 domain, in which two cysteines are spaced 4 amino acids apart with the first cysteine being the last amino acid of the COL1 domain. Such an arrangement is characteristic of the FACIT group of collagens (21–25). The size of the COL1 domain (106 amino acids) is also very similar to that of the COL1 domains of the FACIT members. For example, the COL1 domains of the $\alpha 1(\text{IX})$ and $\alpha 1(\text{XII})$ collagens are 115 and 103 amino acids long, respectively. In addition, each of these COL1 domains contains two Gly-Xaa-Xaa imperfections that are present at similar positions. Furthermore, sequence comparisons reveal that the N-terminal NC11 domain exhibits sequence similarities to the NC domains of two FACIT members—i.e., $\alpha 1(\text{IX})$ and $\alpha 1(\text{XII})$ collagens—and of the fibrillar $\alpha 1(\text{XI})$ collagen chain (Fig. 4). Specifically, a 250-amino acid segment of the NC11 domain shares 27.2% identity with the NC4 domain of the human $\alpha 1(\text{IX})$ collagen (25), 17.6% identity with a segment in the C terminus of the NC3 domain from the chicken $\alpha 1(\text{XII})$ collagen (23, 24), and 19.2% identity with the amino propeptide of the human

Table 1. Sizes (amino acids) of COL and NC domains of $\alpha 1(\text{XVI})$ collagen

	No. of amino acids										
	1	2	3	4	5	6	7	8	9	10	11
NC	26	39	23	34	11	15	21	17	15	14	312
COL	106	422	15	52	138	71	59	34	131	27	—

```

COL16A1 ANTGACQPPSQEGLKLEHSSSLPANVTGFNLIHRLSLM-KKSAIKKIRN
COL9A1 RPRFPVNSNSNGGNELCPKIRIGQDDLPGFDLISQFQV-DKAASRAAIQR
COL11A1 T-TLALTFLLQAREVR-GAAPVDVLKALDFHNSPEGIS--KTTGFCNTRK
COL12A1 IQDNLVTFVCEATATSTCPLIYLEGTYSPGKMLSEYNLTERKHFASVQGVGS

COL16A1 P-KGPLIL----RLGA-APVTPQTRRVFPRG-LPEEPALVTLTLKKKHTH
COL9A1 VV-GSATLQVAYKLGNNVDFRIPTRNLYPSG-LPEEYSFLTTFRMTGSTL
COL11A1 NSKSGSDT---AYRVSQKQLSAPTKQLFPGGTFPEFDESILFTVVKPKKGIQ
COL12A1 LESGSPSYVAYRLHKNAFVSQIETREIHPHG-LPQAYTIIMLFRLLPESP

COL16A1 QKTWYLFQVTDANGYFOISLEVNSQERSLELRAQGD-GDFVSCIEP---
COL9A1 KKNWNIWQIDSSGKEQVGIKINGQTSQVVSFKGLD-GSLQTAASRAISN--
COL11A1 S---FLI-SIYNEHQIQIGVEVGRSPVFL-FEDHTGKPAPEDEYPLFR--
COL12A1 SEFFAIWQITDRDYKPGVGVVLDPSGKVLGSPFNKDR--GEVQTVTFDNDI

COL16A1 VPOLFDRRWHKMLMSVAGRVASVHVDCSSASSOPEG---PRRPMRPVGHVF
COL9A1 LSSLFDSQWKKIMIGVERSSATLFDVDCNRIEISLPK--PRGIDIDGFAV
COL11A1 VN-IADGKWRVAISVEKVTMTIYDCKKTKTKPLDRSERAIIVDNGITV
COL12A1 VKKIYFGSFKVHVI VVTVSSNVKIYIDCSELELKKPK--EAGNITDDGYEI

COL16A1 LG-LDAEQGKPVSFDFLQVQVHIYCDPELVLEEGCEILPAG-CEPETSKAR
COL9A1 LKGLADNPQVSVPELQWNLHCDPLRPRRETCHL-LPARITFSQTTDER
COL11A1 FTGTRLLDEEVFEG-DIQOFLITGDKPKAAY-DYC-EHYSPPD-CDSSAPKAA
COL12A1 LGKLLKGRDRSATLETQNFIVCSVFWTSRDRCCD-LPSMRDEAKPALP
    
```

FIG. 4. Alignment of amino acid sequences of the NC11 domain from the $\alpha 1(\text{XVI})$ collagen (position 1) with the NC4 domain of the human $\alpha 1(\text{IX})$ collagen, amino propeptide of the human $\alpha 1(\text{XI})$ collagen, and the NC3 domain of the chicken $\alpha 1(\text{XII})$ collagen. Shaded regions reflect sequences conserved between the $\alpha 1(\text{XVI})$ collagen and the other collagens. Gaps (—) are introduced to increase identity.

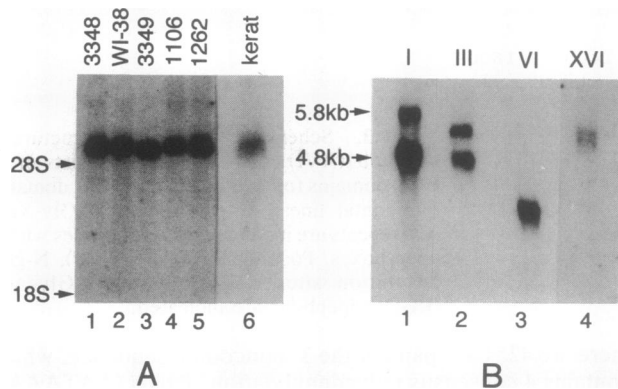


FIG. 5. Northern blot hybridizations of RNAs from cultured human cells with the 4.0-kb cDNA probe F352 (A) and comparison of hybridizations using fibroblast mRNA with cDNA probes for $\alpha 1(I)$, $\alpha 1(III)$, $\alpha 1(VI)$, and $\alpha 1(XVI)$ (B). (A) Each lane contained 2 μ g of the poly(A)⁺ RNA from skin fibroblast 3348 (lane 1), lung fibroblast WI-38 (lane 2), skin fibroblast 3349 (lane 3), fetal skin fibroblast 1106 (lane 4), fetal skin fibroblast 1262 (lane 5), and 20 μ g of total RNA from epidermal keratinocytes (lane 6). Positions of 28S and 18S rRNAs are indicated. (B) One microgram each of poly(A)⁺ RNA from skin fibroblast 3349 was hybridized with cDNA probes for collagen $\alpha 1(I)$ (lane 1), $\alpha 1(III)$ (lane 2), $\alpha 1(VI)$ (lane 3), and $\alpha 1(XVI)$ (lane 4). Autoradiography of lane 4 was ≈ 5 times longer than that of lanes 1–3. The $\alpha 1(I)$ collagen transcripts are 5.8 and 4.8 kb.

$\alpha 1(XI)$ collagen (26). However, only 2 of the 7 cysteines of the NC11 segment are invariant in the other NC domains, indicating that all these NC domains are different in their folding patterns. The predicted amino acid sequence contains three potential N-glycosylation sites (Asn-Xaa-Thr/Ser) and three Arg-Gly-Asp sequences.

Analysis of mRNA Expression. By Northern blot analyses, the cDNA hybridized to a 5.5-kb mRNA species (Fig. 5A). The mRNA was expressed at approximately the same level in all dermal and lung fibroblast cell strains examined. It was also readily detectable in epidermal keratinocytes. However, in fibroblasts the hybridization signal was considerably lower than that obtained by using cDNA probes for types I, III, and VI collagens (Fig. 5B).

Chromosomal Assignment. A total of 100 human metaphase cells were analyzed after *in situ* hybridization of these cells to [³H]cDNA F352 (specific activity, 4×10^7 dpm/ μ g). Of the 246 silver grains associated with chromosomes, 101 grains (41%) were located on chromosome 1. Two representative partial human metaphases with silver grains on G-banded chromosome 1 are shown in Fig. 6A. The chromosomes were subsequently identified by R-banding (Fig. 6A Lower). Approximately 80% of labeled sites on chromosome 1 were found in band p34–35 (Fig. 6B). No other chromosome showed significant labeling in these experiments.

DISCUSSION

The cDNA clones reported here encode a collagenous polypeptide beginning with a hydrophobic signal peptide, suggesting that the predicted protein is secreted into the extracellular space. Approximately 65% of the entire molecule is composed of repeating Gly-Xaa-Xaa sequences, which are represented by 10 separate COL domains. The NC domains contain numerous cysteines, which are often found in the sequence motif of Cys-Xaa-Xaa-Cys (where Xaa is any amino acid). The estimated molecular weight of the unmodified chain is 157,692. The structural characteristics of the putative protein are significantly different from those of the 14 distinct

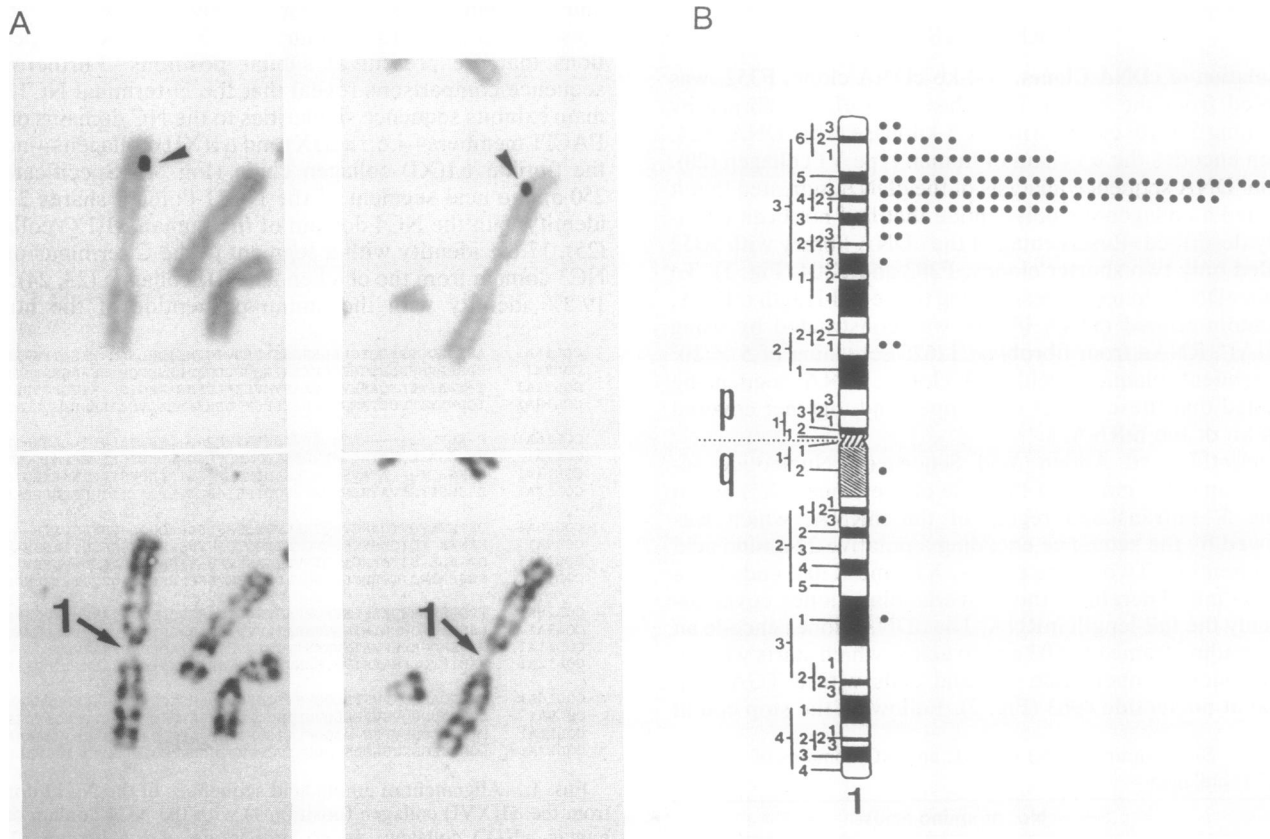


FIG. 6. *In situ* chromosomal mapping of the $\alpha 1(XVI)$ gene to human chromosome 1p34–35. (A) Two partial human metaphase spreads showing the site of hybridization to chromosome 1. (Upper) Presence of silver grains on Giemsa-stained chromosomes after autoradiography is shown. (Lower) Subsequent identification of the chromosome by R-banding. (B) Distribution of silver grains on G-banded chromosome 1.

collagen types described to date and other collagenous sequences reported recently (27, 28). In addition, the cDNA is not homologous to the newly isolated cDNA for type XV collagen (J. C. Myers, personal communication). Since the entire coding region of the mRNA has been characterized, it is highly unlikely that the cDNA clones reported here encode an α chain of a previously described collagen type. Taken together, the data support the conclusion that the cDNA encodes an α chain of a newly discovered collagen type, and we therefore suggest that the cDNA encodes the $\alpha 1$ chain of type XVI collagen.

Many characteristic features of the FACIT group are found in the $\alpha 1(XVI)$ collagen, suggesting that this is a member of the FACIT group. It is noteworthy that the Cys-Xaa-Xaa-Cys motifs in the short NC domains of the $\alpha 1(XVI)$ collagen are not present in other members of the FACIT group. Interestingly, these motifs are found in the short NC domains of several cuticle collagens from *C. elegans*, which resemble the vertebrate FACIT group in having similar COL1 domains (7). It is therefore possible that the $\alpha 1(XVI)$ collagen and the cuticle collagens are evolutionarily related.

In situ chromosomal hybridization demonstrates that the gene encoding the $\alpha 1(XVI)$ collagen is located in the p34-35 region of human chromosome 1. It is of some interest that the $\alpha 1(XI)$ collagen, which shares significant sequence identity with the $\alpha 1(XVI)$ chain, is mapped to the p21 region of chromosome 1 (29).

The cDNA hybridizes to a mRNA of 5.5 kb by Northern blot hybridization analyses. The size is in good agreement with that estimated from the near full-length cDNA. The mRNA is expressed in both dermal and lung fibroblasts but represents only a minor mRNA species in these cells. The message is also present in epidermal keratinocytes, suggesting that type XVI collagen may have a function in the epidermis. In this regard, it is noteworthy that cuticle collagens are expressed in epithelial cells of *C. elegans*. Whether or not these two collagens share similar functions remains to be elucidated.

It has been shown that the FACIT members are localized on the surface of major collagen fibrils and may serve as molecular bridges that are responsible for maintaining the structural integrity of the extracellular matrix. The structural similarities between the $\alpha 1(XVI)$ collagen and the FACIT group raise the intriguing possibility that the $\alpha 1(XVI)$ collagen may serve similar functions. The NC4 domains of both chicken and human $\alpha 1(IX)$ collagen have a basic isoelectric point of ≈ 10 and are thought to be involved in interactions with polyanions in the extracellular matrix (3). The homologous NC11 domain of $\alpha 1(XVI)$ collagen, however, has a neutral isoelectric point of ≈ 7 ; therefore, the functional implication of the sequence identity, if any, remains to be studied. The cloning and characterization of the cDNA presented here will be an important step toward elucidating the biological functions of type XVI collagen.

We thank Dr. Charlene Williams for her critical reading of the manuscript. This research was supported by National Institutes of Health Grants AR38912 and AR38923.

1. Vuorio, E. & de Crombrughe, B. (1990) *Annu. Rev. Biochem.* **39**, 837-872.
2. van der Rest, M. & Garrone, R. (1991) *FASEB J.* **5**, 2814-2823.
3. Shaw, L. M. & Olsen, B. R. (1991) *Trends Biochem. Sci.* **16**, 191-194.
4. Vaughan, L., Mendler, M., Huber, J., Bruckner, P., Winterhalter, K. H., Irwin, M. H. & Mayne, R. (1988) *J. Cell Biol.* **106**, 991-997.
5. Keene, D. R., Lunstrum, G. P., Morris, N. P., Stoddard, D. W. & Burgeson, R. E. (1991) *J. Cell Biol.* **113**, 971-978.
6. Sugrue, S. P., Gordon, M. K., Seyer, G., Dublet, B., van der Rest, M. & Olsen, B. R. (1989) *J. Cell Biol.* **109**, 939-945.
7. Field, C. (1988) *J. Mol. Evol.* **28**, 55-63.
8. Chu, M.-L., Zhang, R.-Z., Pan, T.-C., Stokes, D., Conway, D., Kuo, H.-J., Glanville, R., Mayer, U., Mann, K., Deutzmann, R. & Timpl, R. (1990) *EMBO J.* **9**, 385-393.
9. Chu, M.-L., Gargiulo, V., Williams, C. J. & Ramirez, F. (1985) *J. Biol. Chem.* **260**, 691-694.
10. Rigby, P. W. J., Dieckmann, M., Rhodes, C. & Berg, P. (1977) *J. Mol. Biol.* **113**, 237-251.
11. Sanger, F., Nicklen, S. & Coulson, A. R. (1977) *Proc. Natl. Acad. Sci. USA* **74**, 5463-5467.
12. Pearson, W. R. & Lipman, D. J. (1988) *Proc. Natl. Acad. Sci. USA* **85**, 2444-2448.
13. Chomczynski, P. & Sacchi, N. (1987) *Anal. Biochem.* **162**, 156-159.
14. Thomas, P. S. (1980) *Proc. Natl. Acad. Sci. USA* **77**, 5201-5205.
15. Sambrook, J., Fritsch, E. F. & Maniatis, T. (1989) *Molecular Cloning: A Laboratory Manual* (Cold Spring Harbor Lab., Cold Spring Harbor, NY), 2nd Ed.
16. Chu, M.-L., Myers, J. C., Bernard, M. P., Ding, J.-F. & Ramirez, F. (1982) *Nucleic Acids Res.* **10**, 5925-5934.
17. Chu, M.-L., Weil, D., deWet, W., Bernard, M. P., Sippola, M. & Ramirez, F. (1985) *J. Biol. Chem.* **260**, 4357-4363.
18. Chu, M.-L., Mann, K., Deutzmann, R., Pribula-Conway, D., Hsu-Chen, C. C., Bernard, M. P. & Timpl, R. (1987) *Eur. J. Biochem.* **168**, 309-317.
19. Mattei, M. G., Philip, N., Passage, E., Moisan, J. P., Mandel, J. L. & Mattei, J. F. (1985) *Hum. Genet.* **69**, 268-271.
20. Chu, M.-L., Conway, D., Pan, T.-C., Baldwin, C., Mann, K., Deutzmann, R. & Timpl, R. (1988) *J. Biol. Chem.* **263**, 18601-18606.
21. Ninomiya, Y. & Olsen, B. R. (1984) *Proc. Natl. Acad. Sci. USA* **81**, 3014-3018.
22. Ninomiya, Y., van der Rest, M., Mayne, R., Lozano, G. & Olsen, B. R. (1985) *Biochemistry* **24**, 4223-4229.
23. Gordon, M. K., Gerecke, D. R. & Olsen, B. R. (1987) *Proc. Natl. Acad. Sci. USA* **84**, 6040-6044.
24. Yamagata, M., Yamada, K. M., Yamada, S. S., Shinomura, T., Tanaka, H., Nishida, Y., Obara, M. & Kimata, K. (1991) *J. Cell Biol.* **115**, 209-221.
25. Muragaki, Y., Kimura, T., Ninomiya, Y. & Olsen, B. R. (1990) *Eur. J. Biochem.* **192**, 703-708.
26. Yoshioka, H. & Ramirez, F. (1990) *J. Biol. Chem.* **265**, 6423-6426.
27. Giudice, G. J., Squiguera, H. L., Elias, P. M. & Diaz, L. A. (1991) *J. Clin. Invest.* **87**, 734-738.
28. Marchant, J. K., Linsenmayer, T. F. & Gordon, M. K. (1991) *Proc. Natl. Acad. Sci. USA* **88**, 1560-1564.
29. Henry, I., Bernheim, A., Bernard, M., van der Rest, M., Kimura, T., Jeanpierre, C., Barichard, F., Berger, R., Olsen, B. R., Ramirez, F. & Junien, C. (1988) *Genomics* **3**, 87-90.