



Identification and characterization of the human type II collagen gene (*COL2A1*)

(cartilage collagen gene/mRNA/DNA sequence)

KATHRYN S. E. CHEAH*[†], NEIL G. STOKER*, JANE R. GRIFFIN*, FRANK G. GROSVELD[‡],
AND ELLEN SOLOMON*

*Somatic Cell Genetics Laboratory, Imperial Cancer Research Fund Laboratories, P.O. Box 123, Lincoln's Inn Fields, London WC2A 3PX, United Kingdom; and [‡]Medical Research Council, National Institute for Medical Research, The Ridgeway, London NW7 1AA, United Kingdom

Communicated by Walter Bodmer, December 3, 1984

ABSTRACT The gene contained in the human cosmid clone CosHcoll, previously designated an $\alpha 1(I)$ collagen-like gene, has now been identified. CosHcoll hybridizes strongly to a single 5.9-kilobase mRNA species present only in tissue in which type II collagen is expressed. DNA sequence analysis shows that this clone is highly homologous to the chicken $\alpha 1(II)$ collagen gene. These data together suggest that CosHcoll contains the human $\alpha 1(II)$ collagen gene *COL2A1*. The clone appears to contain the whole gene (30 kilobases in length) and will be extremely useful in the study of cartilage development and for identifying those inherited chondrodystrophies in which defects occur in this gene.

Collagens are major structural components of the extracellular matrix. In vertebrates they form a large family of proteins represented by at least nine distinct types for which a minimum of 17 genes exist to code for their constituent α chains (1-5). Different tissues are characterized by the types and quantity of collagen expressed. The coordinated expression of these different collagen genes is believed to be important in vertebrate development (6), and collagen abnormalities may be involved in a wide range of inherited connective tissue disorders in man (7, 8). To approach these questions, a variety of cDNA and genomic collagen clones from a number of species have been isolated, including the human $\alpha 1(I)$ (9, 10) and $\alpha 2(I)$ genes (11, 12).

We previously reported the isolation of the genomic clone CosHcoll from a human placental cosmid library (13), using the chicken $\alpha 1(I)$ cDNA clone pCg54 as a probe (14). This cosmid clone contains a 36-kilobase (kb) insert and cross-hybridizes with collagen $\alpha 1(I)$ mRNA. However, the amino acid sequence derived from 1 kb of the clone showed only 60-70% homology to chicken and bovine $\alpha 1(I)$ and $\alpha 2(I)$ collagen, and it did not match the human $\alpha 1(III)$ amino acid sequence. Since interspecies protein sequence homologies between collagens of the same type are usually greater than 80%, we concluded that CosHcoll did not code for any of these chains. In the absence of positive identification, we labeled the clone an $\alpha 1(I)$ collagen-like gene.

To establish the identity of this collagen gene, its homologous mRNA was sought and a more extensive nucleotide sequence was obtained. We report here that CosHcoll hybridizes strongly with human fetal cartilage mRNA but not to mRNA from a large number of other sources, suggesting that its expression is cartilage specific. Analysis of the DNA sequence obtained shows that CosHcoll is highly homologous to chicken $\alpha 1(II)$ collagen, which is the major hyaline cartilage collagen. We therefore concluded that CosHcoll probably contains the human $\alpha 1(II)$ collagen gene.

MATERIALS AND METHODS

Enzymes. Restriction endonucleases and DNA-modifying enzymes were purchased from Bethesda Research Laboratories, Boehringer Mannheim, or New England Biolabs.

DNA Preparation, Manipulation, and Sequencing. Standard DNA manipulations were performed as described by Maniatis *et al.* (15). DNA sequencing was carried out as described by Bankier and Barrell (16).

Preparation of Poly(A)⁺ RNA. mRNA was prepared from 10^8 to 5×10^8 cultured cells and from human fetal sterna (16 and 22 weeks) and calvaria (10, 12, and 14 weeks). Cells were homogenized in a solution containing proteinase K (Boehringer Mannheim) at 200 μ g/ml, 20 mM Tris·HCl at pH 7.6, 1 mM EDTA, and 2% sodium dodecyl sulfate. Fetal tissues were frozen in liquid N₂ and ground to a fine powder with a mortar and pestle. The pulverized tissue was then homogenized in the sodium dodecyl sulfate/proteinase K solution described above. Poly(A)⁺ RNA was isolated as described by Cheah *et al.* (17). Usually 1-5% of total RNA was recovered as poly(A)⁺ RNA.

RNA Blot Analyses. mRNAs [1-2 μ g of poly(A)⁺ RNA per gel slot] were denatured with glyoxal, electrophoresed in 0.8% agarose gels, and transferred to filters as described by Thomas (18) except that Pall Biodyne (Santa Monica, CA) A nylon membranes were used. Hybridization of the blots was performed as described (13).

Isolation of Overlapping Cosmid Clones. Overlapping clones from three human cosmid libraries were isolated (19), using the *EcoRI* fragments at the ends of the insert in CosHcoll as probes.

RESULTS

Identification of a Homologous mRNA Species. To establish the identity of CosHcoll, it was necessary to find a homologous mRNA species. mRNAs were prepared from different tissues and cultured cell lines that synthesize characteristic collagen types. These included collagen type I (human fetal calvaria, fibroblast lines) (1), type II (human fetal sterna, rat chondrosarcoma, chicken sterna) (1, 20), type III (human fibroblast lines) (1), type IV (mouse parietal endoderm, human fibrosarcoma line HT1080) (21, 37), type V (human placenta, rhabdomyosarcoma line A204) (1, 22), and type VI (human placenta) (23). The ability of CosHcoll to hybridize to these mRNAs was tested by blot hybridization.

CosHcoll hybridized strongly to mRNA in preparations from only three of the tissues tested: a rat chondrosarcoma, chicken sternal cartilage, and human fetal sterna (Fig. 1). The Swarm rat chondrosarcoma is a transplantable tumor of cartilage origin and has been shown to synthesize pre-

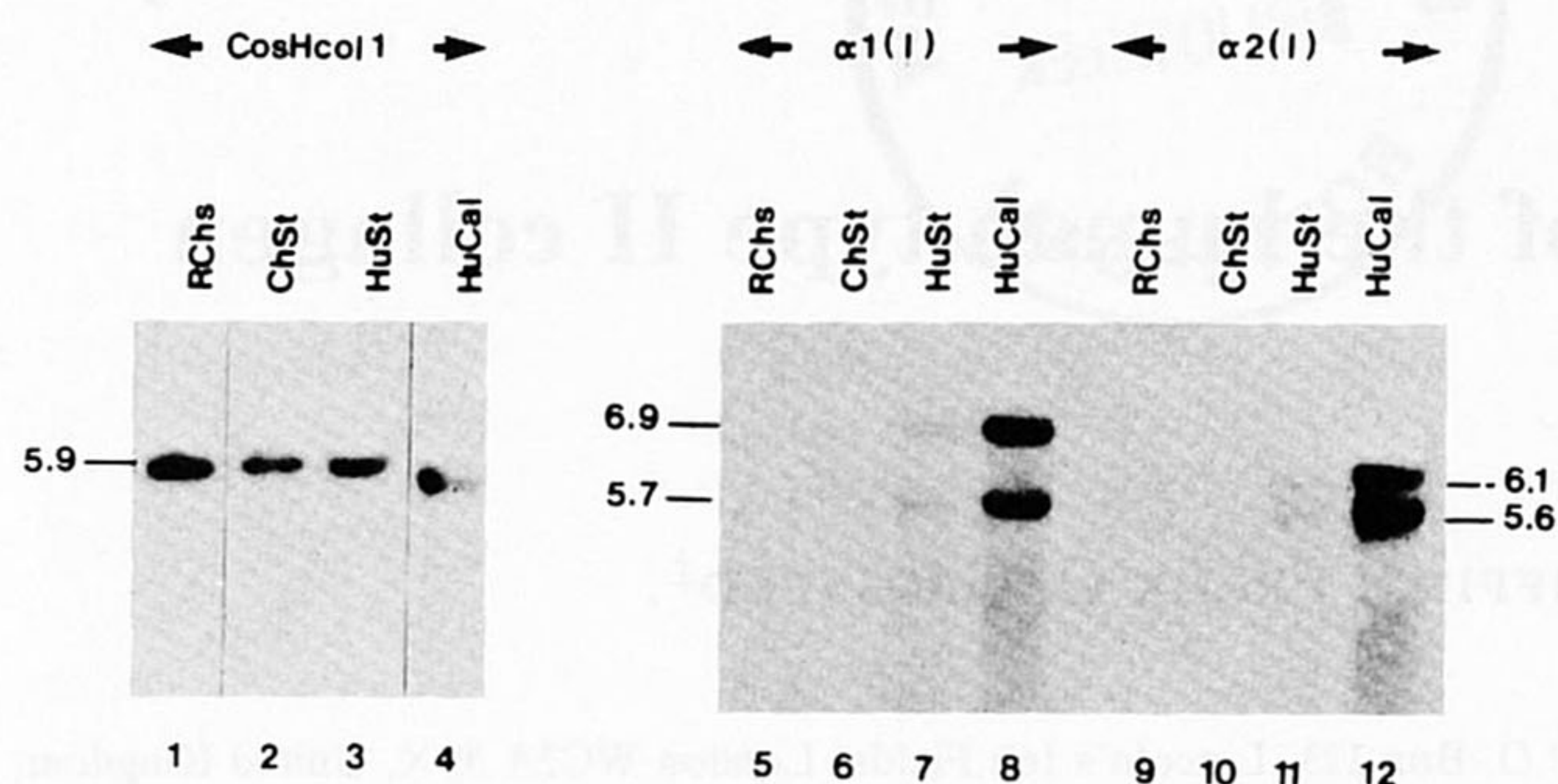


FIG. 1. Hybridization of CosHcoll1 to mRNA preparations. Poly(A)⁺ RNA preparations (1 μ g) from rat chondrosarcoma (RChs), chicken sterna (ChSt), human sterna (and rib ends) (HuSt), and human calvaria (HuCal) were blotted and hybridized with ³²P-labeled nick-translated CosHcoll1, or human α 1(I) or α 2(I) genomic probes. All tracks shown represent overnight exposures except for the hybridization of CosHcoll1 to human calvaria mRNA, which was a 2-day exposure. Sizes are given in kb.

dominantly type II collagen (20), which is the major collagen synthesized by chondrocytes. Chicken sternal cartilage synthesizes mainly type II collagen, as well as other minor types (5, 24, 25). The human fetal sterna used consisted mainly of sternal cartilage with the ends of rib bones attached and would be expected to be synthesizing collagens typical of bone and of cartilage—i.e., types I and II, respectively. CosHcoll1 therefore appeared to hybridize to a cartilage-specific mRNA, possibly type II collagen mRNA.

Fig. 1 shows that CosHcoll1 hybridizes strongly to a distinct 5.9-kb band in mRNA from the cartilaginous tissues (tracks 1–3). The same probe hybridizes less strongly to two different bands of 5.7 and 6.9 kb in mRNA from human fetal calvaria (skull bones) (track 4), which synthesize type I but not type II collagen. These two bands correspond to the sizes of α 1(I) mRNAs of type I collagen. To show that the calvaria were indeed synthesizing type I collagen, and to determine whether type I was also present in the other three tissues, the mRNA preparations were hybridized with human α 1(I) and α 2(I) probes (tracks 5–8 and 9–12, respectively). The human calvaria can be seen to contain large amounts of type I mRNAs (tracks 8 and 12). The human sterna produced small amounts of type I mRNAs (tracks 7 and 11), as was expected from the presence of the ends of the rib bones, but the rat chondrosarcoma and chicken sterna produced virtually no type I collagen. Cross-hybridization of CosHcoll1 to the α 1(I) mRNA was only seen where type I was present in very large amounts, as in human calvaria and chicken tendon (latter not shown). Cross-hybridization to α 2(I) mRNA was not seen at all.

CosHcoll1 did not hybridize with mRNA from other tissues tested, suggesting that it did not code for collagen types III, IV, V, or VI (data not shown).

Nucleotide Sequence Determination. The 3.8-kb *Eco*RI fragment and part of the adjacent 4.3-kb *Eco*RI fragment (see Fig. 4) were sequenced (Fig. 2). Exons were located by comparison with other collagen gene sequences and by following the A-G-G-T splicing rule (38). The nucleotide sequence was combined with that already published (ref. 13; Fig. 2).

The sequenced fragments encode from amino acid 832 to the end of the triple-helical region and the entire C propeptide and extend into the 3' untranslated region. Comparison with other collagen genes shows the CosHcoll1 sequence to be most similar to the chick α 1(II) gene (Table 1), and these are shown aligned in Fig. 2. Where there are differences between the published genomic and cDNA chick α 1(II) sequences (26, 31), the genomic sequence has been used in preference. The derived amino acid sequence from

CosHcoll1 is shown in Fig. 3, aligned with the chicken α 1(II) and human α 1(I) and α 2(I) amino acid sequences. Although the DNA sequence of the chick α 1(II) gene extends only up to exon 4, direct amino acid sequence analyses for exons 5, 6, and 7 show that the high homology continues further (Table 1 and Fig. 3; W. Butler, personal communication). As can be seen from Table 1, the amino acid homologies between CosHcoll1 and the chick α 1(II) gene in exons 1–7 range from 83% to 94% (89% overall), whereas the same exons show only 61–83% (71% overall) and 61–72% (65% overall) homology for the human α 1(I) and α 2(I) chains, respectively. Other published sequences—e.g., chick α 1(III) collagen (30)—all show much lower homology than the chick α 1(II) gene to CosHcoll1 (data not shown). The exon-intron organization of the sequenced region of CosHcoll1 is shown in Fig. 4. The sizes of exons 1–4 are conserved between CosHcoll1 and the chick α 1(II) gene. Intron sizes are different and no significant homology was detected. However, in other collagen genes, both intron and exon sizes have diverged from CosHcoll1, although the locations of introns within the coding sequence have been conserved. We conclude from these results, and from the specific hybridization to human sternal mRNA, that CosHcoll1 codes for α 1(II) collagen.

The 3' Untranslated Region. The sequence of the first 229 bp of the 3' untranslated region of the human α 1(II) collagen gene is shown in Fig. 2. A canonical polyadenylation signal (A-A-T-A-A) is present 189 bp downstream from the stop codon. This, or a similar sequence, is necessary but not sufficient for polyadenylation (36).

Boundaries of the Gene. To determine the extent of the type II gene in CosHcoll1, *Eco*RI fragments from the 5' and 3' regions of the clone were used to screen other human cosmid libraries. Five overlapping clones were isolated, covering a total of 75 kb, of which 12.5 kb was 5' and 25.7 kb was 3' to CosHcoll1. Fragments extending 5' and 3' to CosHcoll1 and fragments from within the clone were tested for hybridization to rat chondrosarcoma mRNA on blots. The results are shown diagrammatically in Fig. 4.

At the 5' end, no hybridization to mRNA was detected with the 9.8-kb *Eco*RI fragment or the 5.9-kb *Eco*RI fragment (which extends 3.2 kb into CosHcoll1). The adjacent 4.8-kb *Eco*RI fragment hybridized to mRNA, as did all the other *Eco*RI fragments in CosHcoll1. The next 3' fragment did not hybridize to mRNA. Since the stop codon and a polyadenylation signal occur within the 3' terminal *Eco*RI fragment of CosHcoll1, and sequences 12.5 kb 5' and 2.5 kb 3' to CosHcoll1 did not hybridize to mRNA, this clone probably contains the complete type II collagen gene. Hybridization of parts of the 4.8-kb *Eco*RI fragment has located

Table 1. Percent sequence homology between CosHcoll1 and other collagen genes (amino acid/nucleotide)

Exon	Chick α 1(II)	Human α 1(I)	Human α 2(I)
10	—/—	67/69	72/63
9	—/—	68/70	53/62
8	—/—	67/56	50/44
7	94/—	83/80	69/58
6	89*/—	61/72	67/72
5	91/—	72/74	61/62
4	91/82	68/68	67/66
3	83/83	71/79	63/73
2	85/84	68/74	62/68
1	94/82	73/70	67/70

Positions where deletion or insertion events have occurred have not been used in this comparison.

*Part of the sequence of this exon is based on amino acid composition and alignment for maximum homology (see Fig. 3 legend).

66TGAACCTG GACGAGAGgt gacagtgag accccctggg gtggccctga ttggggagag gggccctgtg agtctctgtg ctgggtcagc aaggacaagc 100
cccagtcagg gcctcggaga agggggcggc agcgcctggc gacaggcga agccttagta caatgggaag gttgtcgggg agagagacgg gcatagagac 200
caagggtcgc ttctggaagg aggggggaaa cttggtgagg aaactttggc ttcaaatgtg gagtgtgttg ggcagaagag gagaggcctg ggcttctgag 300
aggggctggg ggagcagagg gggaggtgga cagaggacag ctctaggtgc gttcttgttt cactttgtcc ag66AAGCCC C6GT6CT6AT 66CCCCCT6 400
GCASAGAT66 C6CT6CT66A 6TCAAG6gtg gtgtctgtgt tctgtgtgtg cagtgggttg gggaggacat tgcctcgggc ctgacaggtc agctgggggt 500
ggcaggttg aacaagtctc atctcagcct agaaggacct tctgttctct tctcttctgg aacattcttc tctgagcctg agacctctct cctgacag66 600
TAATCGT6GT 6AAACC66T6 CTGT666A6C TCCT66AACC CCT666CCCC CT66CTCCCC T66CCCC6CT 66TCCAAC6T 6CAAGCAAG6 6GACAG66A 700
6AAGCTgtaa gtatcctgga attcagtaaa agccgccttc ccctgcgcgg tggggctgag gcagtccttg ggtttccgca gtctctggac taaggagcag 800
tggcctcaga tgcagagagg gccccacct gtcctggctt ttctctgacg ctgcgcctac tctctctca g66T6CACA 66CCCCAT66 6ACCC6CA66 900
ACCA6CT66A 6CCCC666AA TCCA6gtgag tatccaagt tctctgactg agtccccacc agggatagcg tgggagggca gccagcctcc aggtggttcc 1000
tggcctccag ccctgtgttt ccggggattc ctgagcttgg gtgggacagg agggggctcc tgcctctggc ctgacctgac tcaatcggtg tctgtctgt 1100
tcccag66T6 CTCA666CCC CA6A66T6AC AA66A6A666 CT66A6A66C T66C6A6A6A 66CCT6A666 6ACACC6T66 CTTACT66T CT6CA666T6 1200
T6CCC666CC TCCTgtgagt gtcactgctt gcgtgggact tcccagggcc tctctgccca cagagcccac ttgagctccc tgtctgcca ggacagcttg 1300
ggatcacctt aagcagtttc taggatttcc tcagggttgg agggaggagg aagtggaaag ggaatggggc tgggacataa agctgttccc ccagctccca 1400
gaatatagat agatagtct gtctgacgg tggcctttt cctcttctt ctacacag66 TCCTTCT66A 6ACCA66T6 CTTCT66TCC T6T66TCC 1500
TCT66CCCTA 6Agtaaagtga catggagttg gaagatggag gggggccttc agagagtggt ggcctgtgtt cccatgggga gggaaatgct gctgcttctg 1600
gggaagctgt gggctcaggg gtcctcactc agtaatgggg gcaggacttg ctcatgtgct tatggccaga aaagcgcctg aggccacaat ggctgtaaga 1700
caaacatgaa tcagcctctc gctgtcagac agaacagcat ttacaaaaga ggagcttagg agggtaggca agccatggag ctatctgct ggttcttggc 1800
caaatagaga ccaacttagg gttccatgac tgagcatgtg aagaactggg ggcggagtggt ctgggtctat caggacagcc acctaccag ccccagcagc 1900
tccccagcct tccctgtgtt gaccactctt tctctcagc ctctctctct tgcag66T6C TCCT66CCCC 6T66T6CCT CT66CA66A T66T6CTAAT 2000
66AATCCCT6 6CCCCAT66T 6CCTCCT66T 6CCCC666AC 6ATCA666CA AAC666CCT 6CTgtaaagt tctctgactc tctctctgtg tggaggtgct 2100
cctaccatcc gggaggttg agctcttttt tgcctcaggg ctcttttagg gcctcagcct gcagctaaac gtgatggcat cctttatctt gaggctcct 2200
cagaggtcac agggcccatg atcagtgctg ggaactgaa gagaagggtc aaggaagaaa tagacatggt gctgtgtgtt ccttggctct cgcctgctac 2300
acctccgcc caccatggg gctgggaaga gggacactct agtacattct agcaaatggg gatggacatg gaggggcact ttcacacaat cctggctgat 2400
ctctctgttt cctgtctgag 66TCTCCT6 6AATCCT66 ACCCCCT66T CCTCC666T 6CCCC666CC T66CAT66C AT6TCC66CT T66T66CTT 2500
.....C..C.. T..T..C... ..T..C.C...A. C..... ..T..T.AC.
A66CC6A6A 6A6A6666C 6C6ACCCCT 6C6ATACAT6 C666C6ACC 666C66C66 T66CCT6A6A 6AC6AT6AC6 6C6A66T66A T66C6ACATC 2600
G..T.A..C6A. C.6C..... A...A...6 ...6..... A..6...C.6C.... T6..... C.....C...
AAGTCCCTCA ACAAC6A6AT 6A6A66C6AT 6C6A66C66C 6666CCTCC6 CA6A66CCT 6CTC66CCT 6C6A66CCT 6AAACT66C 6ACCC6A6T 2700
..A..... ..T..... AA 6..... ..CA.6... ..C.C...A. C..... ..T..C....
66A6A6AT6g taagcttggg gaacaggatc ccctgcccgg ggaagcaggg agtcatcctt taggcctagc agcaaggagg gagatgccc ctagtacagg 2800
.....C.
gcagagctgg gcctggaagt ttccgccaga gggttcctct cttatttca acgagagaag ctgacgctt gcccctgtc ctgcatggc tacctggcgg 2900
aggtgacctc aggtggact ccatccacca gctgggact gcttctgctc tctttgcatg tgttcttct tagggctgga cttagctcat gcagatctcc 3000
ctgcccctgc atctcccag gtcctccctc ttcaggcca catgtgaacc tcatcccttg tccctgtagg cctctctgtc tctttcagtc aggcctgggt 3100
ctctcaagct tttgtgtctg tgcctgtctg agcccccatg ggtgtgctct cttccccctg cag6A6ACTA CT66ATT6AC CCCAACC6A6 6CT66C6CTT 3200
.....T..6.....6.
66AC6CCAT6 A66T6TTTCT 6CAACAT66A 6ACT66C6A6 ACTT66CTCT ACCCCAATCC 66CAAC6GT CCCA6A6A6A ACT66T66A6 6A6CA6A6C 3300
.....C ..A..A...A..... ..C..... ..G.CC.. 6A6C.6.A.C6.....CC6
A66A6A6A6A AACACATCT6 6TTT66A6AA ACCATCAAT6 6T66CTTCCA Tgtgagtagc tgggtgcccct agatgatgag cagagatggc tctcaaac 3400
..A..C.... .6...6...C...6C. .C..T.... C
ctttcttttc tttctccctg gaagcttita gcaacttccc catattttcc tccagtttcc tgttgggctt gagaggaggg aaaagtaggg aaaagtattt 3500
tttccccagc tggaggtggg aaaagaggtc ctctgagctt gctccactcc tggaaagaaa aatgtccaac tagctccctg ctgccccagc acctttagg 3600
tcttgaacc atgaactctt ggcagcccct acagcccctg gtcctattga atgcaagctc ccaggcctca cactgcccct ctctgcccc acaTTC6C 3700
.....
TAT6A6AT6 ACAATCT66C TCCCAACACT 6CCAAC6TCC 6AT6ACCTT CCTAC6CCT6 CT6TCC6C6 6A66CTCCCA 6A6ATCACC TACCACT6CA 3800
..C..C.... .6..C...T. C.....C6.A...6.....CC. .6..... ..6.....
A6A6A6CAT T6CCTATCT6 6AC6A66C6A 6T66CAACT CA6A666C CT6CTCATCC 6666CTCAA T6C6T66A6 ATCC666C6A 6666CAAT6 3900
.....C.....CA..6.A6A .6.....6..... A..A... A.C..... ..A..A... C..... ..A..A..C.C..
C66TTC6C6 TACT66CCC T6A66A6T66 CT6C6C6gtg agtggggctg ccagagagaa gagctgctg tgcctcaact gcctggagca gggctgaggg 4000
.....C6C.T.T ..6....C... ..
ttggcccgcg gcagctgtca ggtcctaaag tgacaggatc atcagagga tgagttttag ggtcatgtag agaagatagg ctgagtgaca ggtgagagag 4100
aggcacatat cattccatct tctccattcc cctggctcag ggaacaaaaa ccctacctg aaccagtgca ctactgtaga agtgttctc caatgtgtac 4200
aggtggaaga agcggtcaca ggttgggagc tcaactgtgg gagtgggaa ggaagggaa ggcaggttg agaagggccc tgcctgtaag gataggagt 4300
gaagtggaga ggcctttggc aagccaagaa gaggctcag gacccccctc agtgtgttcc aacctgtggt gctctgatgc tgcaggttt gttcagttt 4400
gggcttctgg gcagctggaa ctgggtagca aggcactac tgaacagagc ctctctctt tttctccct agAAACATAC C66TA6T66 66CA6ACT6 4500
.....C.. T..C..A...6.
TTATC6A6TA CC66T6C6A6 A6A6CCTC6 6CCTCCCAT CATT6ACATT 6CACC6T66 ACAT66A66 6CC6A6C6A 6AATTC66T6 T66CAT66 4600
..6..... ..6... ..6.....6.6..... T6.A..T... ..T... ..T..C.. 66...T... ..6..T..C.T..T..
6CC66T6C6 TTCTT6TAAA AACCT6AACC 6A6A6A6A6A ACAATCC6T 6CA6A6C6A 66A6C6A66 TACTTCCAA TCTC6T6C TCTA66ACTC 4700
C..A.....***
T6CACT6AAT 66CT6ACT6 ACCT6AT6C CATT6ATCC ACCCTC6C 66T66ACT TTTCTCCCT CTCTTCTAA 6A6ACT6AA CT666C6A6C 4800
T6CAAAATAA AATCT66T6 TTCTATTAT TTATT6CTT CCT6T
***** *
POLY A

FIG. 2. Nucleotide sequence from CosHcoll1. The sequence of the 3.8-kb *EcoRI* fragment and part of the 4.3-kb *EcoRI* fragment (see Fig. 4) was determined. This was combined with the sequence previously published (13), which extends 720 base pairs (bp) into the 9.3-kb *EcoRI* fragment. A few errors in the earlier sequence have been corrected. Of the sequence not previously published, 93% of the protein-coding and 3' untranslated regions and 70% of the intron sequences were determined on both strands. Upper-case letters, exons; lower-case letters, introns. Exons 1-4 are compared with the chicken $\alpha 1(\text{II})$ sequence (26). Only bases that differ from the CosHcoll1 sequence are shown. The termination codon is marked at position 4617, and a canonical poly(A) addition signal is marked at 4806.

all of the mRNA homology within this fragment to a 1.3-kb region beginning 1 kb from the 5' end (data not shown). This information, combined with the location of the putative polyadenylation signal, provides us with an estimate of the gene length of 30 kb.

DISCUSSION

We have identified the human cosmid clone CosHcoll1. Strong hybridization to a 5.9-kb cartilage-specific mRNA and comparison with the chick $\alpha 1(\text{II})$ collagen gene sequence suggest that CosHcoll1 contains the human $\alpha 1(\text{II})$ collagen gene (*COL2A1*), of which we have sequenced 10 exons at the 3' end. This represents approximately 15% of the estimated

gene length and over 30% of the protein-encoding sequence. mRNA hybridization and DNA sequence data together provide evidence that CosHcoll1 may contain the entire human type II collagen gene and that the gene is approximately 30 kb in length.

The isolation of a genomic human $\alpha 1(\text{II})$ collagen clone has recently been reported elsewhere (39), and the published sequence of the 3' end of exon 4 and of the small fragment of adjacent intron matches exactly with CosHcoll1. This suggests that we have cloned the same gene even though no homologous mRNA was described in that report. Furthermore, a 540-bp cDNA clone has recently been isolated from human fetal cartilage (E. Vuorio, personal communication)

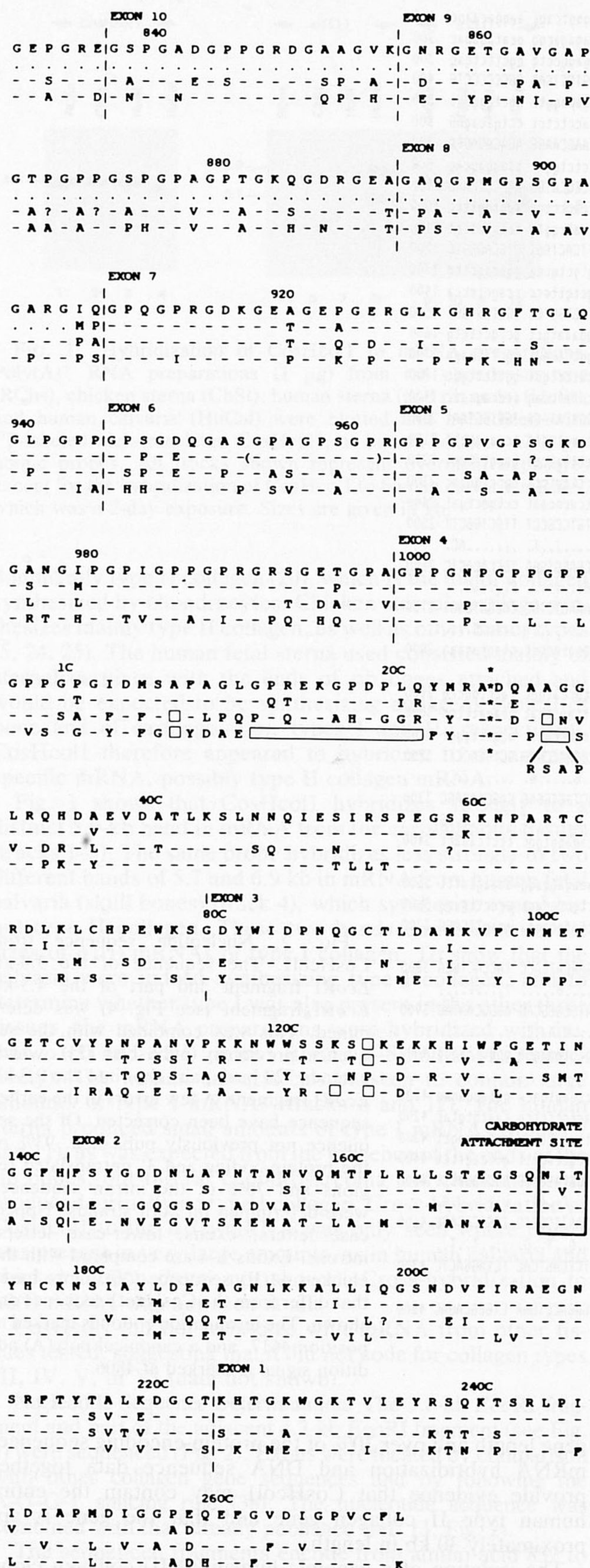


FIG. 3. Amino acid sequence encoded by CosHcoll1. The amino acid sequence encoded by CosHcoll1 was deduced from the nucleotide sequence, and is shown here, compared with other sequences. The standard one-letter code (27) is used. Line 1, CosHcoll1-derived amino acid sequence; line 2, chicken $\alpha 1(\text{II})$ amino acid sequence (26); line 3, human $\alpha 1(\text{I})$ amino acid sequence (28); line 4, human $\alpha 2(\text{I})$ amino acid sequence (29). A dash indicates the presence of the

that is identical in DNA sequence with CosHcoll1, extending from exon 1 to exon 4. The poly(A)⁺ mRNA from which the cDNA clone was derived was shown to program the synthesis of $\alpha 1(\text{II})$ collagen *in vitro* (40). This strongly supports the idea that CosHcoll1 does not carry a pseudogene.

It remains possible that CosHcoll1 carries an $\alpha 1(\text{II})$ -related gene. For example, the minor cartilage collagen chain 3α is highly homologous to $\alpha 1(\text{II})$ collagen (25, 41, 42) and may or may not be genetically distinct. However, no evidence of other sequences homologous to CosHcoll1 has been found in Southern hybridizations, under conditions in which cross-hybridization with the $\alpha 1(\text{I})$ gene was visible (43), and copy number estimates are consistent with only one copy of the gene per haploid genome (R. Dalglish, personal communication). It has been claimed that the 3α collagen chain differs from $\alpha 1(\text{II})$ collagen in that it has a much larger peptide in place of the $\alpha 1(\text{II})$ cyanogen bromide peptide CB9,7 (42). The sequence presented in this paper covers the whole of the region encoding CB9,7 and agrees with the human $\alpha 1(\text{II})$ cyanogen bromide map (44). Nevertheless, absolute proof of this gene's identity will require a comparison of the amino acid sequence of human type II collagen with that derived from the DNA sequence.

The $\alpha 1(\text{II})$ gene has been assigned to chromosome 12 (43, 45). The gene is therefore not linked to the $\alpha 1(\text{I})$ or $\alpha 2(\text{I})$ collagen genes, which map to chromosomes 17 and 7, respectively (46-49), or to the $\alpha 1(\text{III})$ or $\alpha 1(\text{IV})$ collagen genes, which map to chromosomes 2 and 13, respectively (43).

The isolation of the $\alpha 1(\text{II})$ collagen gene is a major step towards the identification of those connective tissue disorders for which an abnormality in this gene is the primary defect. CosHcoll1 should prove to be particularly useful in this respect because it appears to carry the entire gene. Several polymorphisms with high allele frequencies have been identified in this gene (50-52) and are being used for linkage analyses in families with some of these disorders.

We thank the following for their generous gifts: David Rowe and Raymond Dalglish for human $\alpha 1(\text{I})$ and $\alpha 2(\text{I})$ DNA probes, respectively; Markku Kurkinen for parietal endoderm and HT1080 mRNAs; Claudio Schneider for human placental mRNA; Roger Mason, Lance Liotta, and Bryan Sykes for Swarm rat chondrosarcoma tissue, A204 cell line, and human fetal tissue, respectively; and Linda Sandell, Bill Butler, and Eero Vuorio for chicken $\alpha 1(\text{II})$ DNA sequence, chicken $\alpha 1(\text{II})$ amino acid sequence, and human $\alpha 1(\text{II})$ cDNA sequence, respectively. We are grateful to Frances Benham, Elizabeth Weiss, Richard Flavell, Adrian Kelly, Toby Gibson, and Mike Owen for helpful discussions and technical advice. This work was supported in part by a grant from Hong Kong University and the Croucher Foundation, Hong Kong.

1. Bornstein, P. & Sage, H. (1980) *Annu. Rev. Biochem.* **49**, 957-1103.
2. Furthmayr, H., Wiedemann, H., Timpl, R., Odermatt, E. & Engel, J. (1983) *Biochem. J.* **211**, 303-311.
3. Bentz, H., Morris, N. P., Murray, L. W., Sakai, L. Y., Hollister, D. W. & Burgeson, R. E. (1983) *Proc. Natl. Acad. Sci. USA* **80**, 3168-3172.
4. Sage, H., Trueb, B. & Bornstein, P. (1983) *J. Biol. Chem.* **258**, 13391-13401.
5. Ninomiya, Y. & Olsen, B. R. (1984) *Proc. Natl. Acad. Sci. USA* **81**, 3014-3018.

same residue as in the CosHcoll1-derived sequence. Numbers refer to the position in the helix, or in the carboxyl nonhelical domain, and are based on the $\alpha 1(\text{I})$ sequence; the numbers begin above the residue to which they refer. All sequences were derived from nucleotide sequences except for the chicken $\alpha 1(\text{II})$ residues 908-999, which have been determined directly (W. Butler, personal communication). The order of residues 955-962 within this is not known, but the amino acid composition is known, and residues have been aligned for maximum homology. . . . , region not sequenced; ?, a residue not confirmed; *, end of the mature collagen molecule; |, exon boundary; open boxes, deletions/insertions introduced for maximum interchain homology. The carbohydrate attachment site, which lies within a highly conserved region at the nucleotide level (30), is shown.

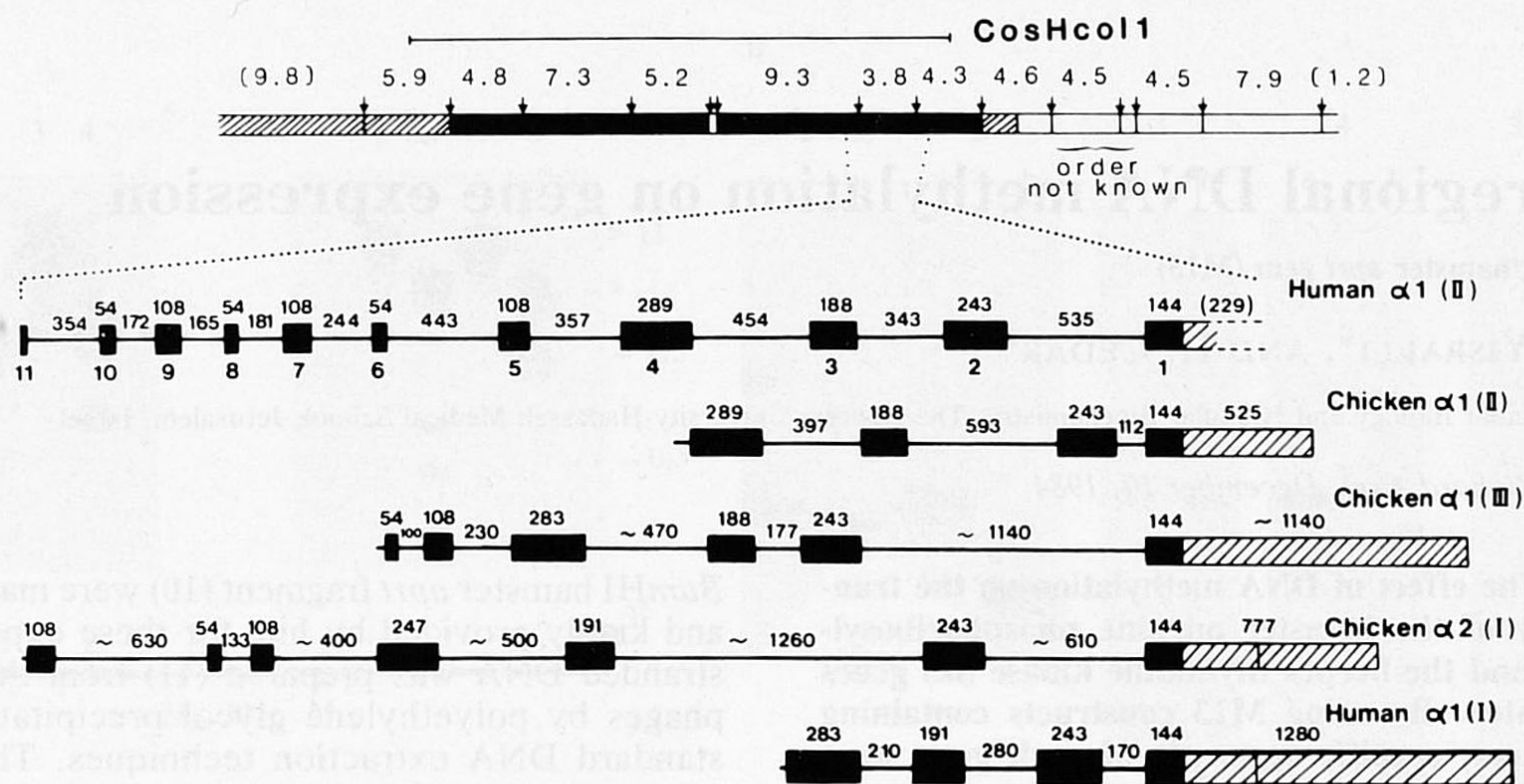


FIG. 4. Organization of the human $\alpha 1(II)$ collagen gene. (Upper) Restriction map showing CosHcol1 within 75 kb of genomic DNA. Fragment sizes (in kb) and positions are composites from five different cosmid clones overlapping CosHcol1 and extending 5' and 3' to it. The positions of *EcoRI* sites are indicated by the arrows. The order of the 4.5- and 1.1-kb fragments (bracketed) was not determined. Filled boxes, *EcoRI* fragments hybridizing to rat chondrosarcoma mRNA (data not shown); hatched boxes, fragments that do not hybridize to mRNA; open boxes, fragments not tested. Only 2.5 kb of the 4.6-kb *EcoRI* fragment was tested. Repetitive sequences were present in all fragments represented by hatched or empty boxes except for the 0.8-kb (the small fragment in CosHcol1) and 1.2-kb fragments, and also in the 5.2-kb fragment. (Lower) The region sequenced is expanded to show its organization into exons and introns (sizes given in bp) and is compared with the chicken $\alpha 1(II)$ gene (26), the chicken $\alpha 1(III)$ gene (30, 32), the chicken $\alpha 2(I)$ gene (30, 33–35), and the human $\alpha 1(I)$ gene (10). Filled boxes indicate protein-coding regions, while hatched boxes indicate 3' untranslated regions. The end of this region has not been determined for CosHcol1. Vertical lines within the 3' untranslated regions indicate the presence of alternative polyadenylation sites.

6. Adamson, E. D. (1982) in *Collagen in Health and Disease*, eds. Weiss, J. B. & Jayson, M. I. V. (Churchill-Livingstone, London), pp. 218–243.
7. McKusick, V. A. (1972) *Heritable Disorders of Connective Tissue* (Mosby, St. Louis, MO), 4th Ed.
8. Hollister, D. W., Byers, P. H. & Holbrook, K. A. (1982) *Adv. Hum. Genet.* **12**, 1–87.
9. Chu, M.-L., Myers, J. C., Bernard, M. P., Ding, J.-F. & Ramirez, F. (1982) *Nucleic Acids Res.* **10**, 5925–5934.
10. Chu, M.-L., de Wet, W., Bernard, M., Ding, J.-F., Morabito, M., Myers, J., Williams, C. & Ramirez, F. (1984) *Nature (London)* **310**, 337–340.
11. Myers, J. C., Chu, M.-L., Faro, S. H., Clark, W. J., Prockop, D. J. & Ramirez, F. (1981) *Proc. Natl. Acad. Sci. USA* **78**, 3516–3520.
12. Myers, J. C., Dickson, L. A., de Wet, W. J., Bernard, M. P., Chu, M.-L., Di Liberto, M., Pepe, G., Sangiorgi, F. O. & Ramirez, F. (1983) *J. Biol. Chem.* **258**, 10128–10135.
13. Weiss, E. H., Cheah, K. S. E., Grosveld, F. G., Dahl, H. H. M., Solomon, E. & Flavell, R. A. (1982) *Nucleic Acids Res.* **10**, 1981–1994.
14. Lehrach, H., Frischauf, A. M., Hanahan, D., Wozney, J., Fuller, F. & Boedtker, H. (1979) *Biochemistry* **18**, 3146–3152.
15. Maniatis, T., Fritsch, E. F. & Sambrook, J. (1982) *Molecular Cloning: A Laboratory Manual* (Cold Spring Harbor Laboratory, Cold Spring Harbor, NY).
16. Bankier, A. T. & Barrell, B. G. (1983) *Techniques in the Life Sciences* (Elsevier, Limerick, Ireland), Vol. B5, pp. 1–34.
17. Cheah, K. S. E., Grant, M. E. & Jackson, D. S. (1979) *Biochem. Biophys. Res. Commun.* **91**, 1025–1031.
18. Thomas, P. S. (1980) *Proc. Natl. Acad. Sci. USA* **77**, 5201–5205.
19. Grosveld, F. G., Dahl, H. M., de Boer, E. & Flavell, R. A. (1981) *Gene* **13**, 227–237.
20. Smith, B. D., Martin, G. R., Miller, E. J., Dorfman, A. & Swarm, R. (1975) *Arch. Biochem. Biophys.* **166**, 181–186.
21. Kurkinen, M., Barlow, D. P., Helfman, D. M., Williams, J. G. & Hogan, B. L. M. (1983) *Nucleic Acids Res.* **11**, 6199–6209.
22. Alitalo, K., Myllylä, R., Pritzl, P., Vaheri, A. & Bornstein, P. (1982) *J. Biol. Chem.* **257**, 9016–9024.
23. Odermatt, E., Ristell, J., van Delden, V. & Timpl, R. (1983) *Biochem. J.* **211**, 295–302.
24. Miller, E. J. (1971) *Biochemistry* **10**, 1652–1659.
25. Burgeson, R. E. & Hollister, D. W. (1979) *Biochem. Biophys. Res. Commun.* **87**, 1124–1131.
26. Sandell, L. J., Prentice, H. L., Kravis, D. & Upholt, W. B. (1984) *J. Biol. Chem.* **259**, 7826–7834.
27. IUPAC-IUB Commission on Biochemical Nomenclature (1968) *Eur. J. Biochem.* **5**, 151–153.
28. Bernard, M. P., Chu, M.-L., Myers, J. C., Ramirez, F., Eikenberry, E. F. & Prockop, D. J. (1983) *Biochemistry* **22**, 5213–5223.
29. Bernard, M. P., Myers, J. C., Chu, M.-L., Ramirez, F. & Eikenberry, E. F. (1983) *Biochemistry* **22**, 1139–1145.
30. Yamada, Y., Kuhn, K. & de Crombrughe, B. (1983) *Nucleic Acids Res.* **11**, 2733–2744.
31. Ninomiya, Y., Showalter, A. M., van der Rest, M., Seidah, N. G., Chretien, M. & Olsen, B. R. (1984) *Biochemistry* **23**, 617–624.
32. Yamada, Y., Liau, G., Mudryj, M., Obici, S. & de Crombrughe, B. (1984) *Nature (London)* **310**, 333–337.
33. Dickson, L. A., Ninomiya, Y., Bernard, M. P., Pesciotta, D. M., Parsons, J., Green, G., Eikenberry, E. F., de Crombrughe, B., Vogeli, G., Pastan, I., Fietzek, P. P. & Olsen, B. R. (1981) *J. Biol. Chem.* **256**, 8407–8415.
34. Wozney, J., Hanahan, D., Morimoto, R., Boedtker, H. & Doty, P. (1981) *Proc. Natl. Acad. Sci. USA* **78**, 712–716.
35. Aho, S., Tate, V. & Boedtker, H. (1983) *Nucleic Acids Res.* **11**, 5443–5450.
36. Fitzgerald, M. & Shenk, T. (1981) *Cell* **24**, 251–260.
37. Pihlajaniemi, T., Myllylä, R., Alitalo, K., Vaheri, A. & Kivirikko, K. I. (1981) *Biochemistry* **20**, 7409–7415.
38. Breathnach, R. & Chambon, P. (1981) *Annu. Rev. Biochem.* **50**, 349–383.
39. Strom, C. M. & Upholt, W. B. (1984) *Nucleic Acids Res.* **12**, 1025–1038.
40. Vuorio, E., Elima, K., Pulkkinen, J. & Viitanen, A.-M. (1984) *FEBS Lett.* **174**, 238–242.
41. Furuto, D. K. & Miller, E. J. (1983) *Arch. Biochem. Biophys.* **226**, 604–611.
42. Eyre, D. R., Wu, J.-J. & Woolley, D. E. (1984) *Biochem. Biophys. Res. Commun.* **118**, 724–729.
43. Solomon, E., Hiorns, L. R., Spurr, N., Kurkinen, M., Barlow, D., Hogan, B. L. M. & Dagleish, R. (1985) *Proc. Natl. Acad. Sci. USA* **82**, in press.
44. Miller, E. J. (1984) in *Extracellular Matrix Biochemistry*, eds. Piez, K. A. & Reddi, A. H. (Elsevier, New York), pp. 41–81.
45. Solomon, E., Hiorns, L., Cheah, K. S. E., Parkar, M., Weiss, E. & Flavell, R. A. (1984) *Cytogenet. Cell Genet.* **37**, 588.
46. Huerre, C., Junien, C., Weil, D., Chu, M.-L., Morabito, M., Van Cong, N., Myers, J. C., Foubert, C., Gross, M.-S., Prockop, D. J., Bone, A., Kaplan, J.-C., De la Chapelle, A. & Ramirez, F. (1982) *Proc. Natl. Acad. Sci. USA* **79**, 6627–6630.
47. Solomon, E., Hiorns, L., Sheer, D. & Rowe, D. (1984b) *Ann. Hum. Genet.* **48**, 39–42.
48. Junien, C., Weil, D., Myers, J. C., Van Cong, N., Chu, M.-L., Foubert, C., Gross, M.-S., Prockop, D. J., Kaplan, J.-C. & Ramirez, F. (1982) *Am. J. Hum. Genet.* **34**, 381–387.
49. Solomon, E., Hiorns, L., Dagleish, R., Tolstoshev, P., Crystal, R. & Sykes, B. (1983) *Cytogenet. Cell Genet.* **35**, 64–66.
50. Driesel, A. J., Schumacher, A. M. & Flavell, R. A. (1982) *Hum. Genet.* **62**, 175–176.
51. Sykes, B. (1983) *Dis. Markers* **1**, 141–146.
52. Sykes, B., Smith, R., Vipond, S., Paterson, C., Cheah, K. & Solomon, E. (1985) *J. Med. Genet.*, in press.