

Molecular cloning and sequence analysis of a chondroitin sulfate proteoglycan cDNA

(extracellular matrix/cell adhesion/glycosylation)

MARIO A. BOURDON, ÅKE OLDBERG, MICHAEL PIERSCHBACHER, AND ERKKI RUOSLAHTI

Cancer Research Center, La Jolla Cancer Research Foundation, 10901 North Torrey Pines Road, La Jolla, CA 92037

Communicated by Phillips W. Robbins, October 19, 1984

ABSTRACT We report the identification and DNA sequence of a chondroitin sulfate proteoglycan core protein cDNA. A cDNA clone, pPG1, was selected from a rat yolk sac tumor poly(A)⁺ RNA-derived cDNA library by using synthetic oligonucleotides predicted from the NH₂-terminal peptide sequence of the mature chondroitin sulfate proteoglycan. The resulting sequence analysis demonstrated that the 874-base-pair pPG1 clone contained the complete coding region of the mature proteoglycan core protein as well as 5' and 3' flanking sequences. The 104 amino acid proteoglycan core protein sequence reveals that the core protein is composed of three regions, the most striking of which is the central 49 amino acid region composed of alternating serine and glycine residues. This region clearly functions as the acceptor site for the attachment of chondroitin sulfate side chains. The serine-glycine repeat region is flanked by a 14 amino acid NH₂-terminal region identical to the NH₂-terminal sequence of the proteoglycan obtained by amino acid sequencing and a 41 amino acid COOH-terminal region. RNA transfer blot hybridizations of poly(A)⁺ mRNA from rat yolk sac tumor cells with nick-translated pPG1 reveal a single mRNA of ≈1300 nucleotides. The possibility of detecting mRNAs and genomic sequences for other proteoglycans with a serine-glycine repeat by using this cDNA clone is discussed.

Proteoglycans composed of a core protein and multiple glycosaminoglycan chains covalently attached to it serve as structural components of various extracellular matrices and basement membranes (reviewed in ref. 1). Proteoglycans are also present at cell surfaces (2-6), interact with other extracellular matrix proteins (7-10), and are thought to be involved in cellular adhesion (11). Relatively little is known about the structure and functions of the various proteoglycans. However, it has become clear that there are several different core proteins for each major type of proteoglycan and that different cell types express different proteoglycans (1, 4, 6, 12, 13). In addition, there may be differences in proteoglycan expression between normal and neoplastic cells (14-16).

Understanding the functions of proteoglycans will require more information about their core protein structure. However, both biochemical and immunological approaches have been exceedingly difficult due to the presence of the abundant and heterogeneous glycosaminoglycan side chains and poor immunogenicity of proteoglycans. Recombinant DNA methods offer a solution to these difficulties by allowing the identification and sequencing of specific proteoglycan mRNAs and genes and, in turn, providing information on amino acid sequence and peptide structure.

We have taken this molecular approach to study a rat yolk sac tumor chondroitin sulfate proteoglycan (17). This proteoglycan interacts with collagen, fibronectin (18), and vi-

tronectin (19) and is capable of inhibiting cell adhesion (20). We describe here the isolation and sequencing of the proteoglycan cDNA. Analysis of the cDNA provides information regarding the extended amino acid sequence and primary structure of a proteoglycan core protein.

MATERIALS AND METHODS

Cell Lines. Rat L2 yolk sac tumor cells (21) and rat hepatoma 7777 cells (22) were cultured in Dulbecco's modified Eagle's medium (Flow Laboratories) supplemented with 10% fetal bovine serum.

Amino Acid Sequencing. Edman degradation of purified proteoglycan was performed with the Applied Biosystems (Foster City, CA) model 470A gas-phase sequencer using the trifluoroacetic acid chemistry provided by the manufacturer. The phenylthiohydantoin amino acids were identified and quantitated by using the Perkin-Elmer series 3B HPLC and ultraviolet detection.

Isolation of Poly(A)⁺ mRNA. RNA was extracted from L2 and 7777 cells by the guanidinium isothiocyanate procedure (23). Poly(A)⁺ RNA was isolated from total RNA preparations by chromatography on oligo(dT)-cellulose (Collaborative Research, Waltham, MA) (24).

Synthesis of cDNA. Double-stranded cDNA was synthesized by using oligo(dT) (P-L Biochemicals) for priming first strand synthesis. Avian myeloblastosis virus reverse transcriptase (Life Sciences, St. Petersburg, FL) was used for synthesis of both strands (25). S1 nuclease-treated cDNA was ultracentrifuged on a 5-20% sucrose gradient and cDNA between 600 and 4000 base pairs (bp) was collected for dC-tailing. Poly(dC)-tailing of the 3' cDNA termini was done by using terminal deoxynucleotidyl transferase (New England Nuclear). The dC-tailed cDNA was annealed to *Pst*I-restricted dG-tailed pBR322 (New England Nuclear) (26) and used to transform *Escherichia coli* K-12 strain MM294 (endoI r_kM_k⁻Blpro⁻) by the procedure of Kushner (27).

Synthesis of Mixed Oligonucleotides. Oligonucleotides were synthesized by the solid-phase phosphotriester method as described (28, 29). The oligonucleotides were purified by HPLC (30) on an RP-C₁₈ reversed phase column (Bio-Rad). Radiolabeled oligonucleotides were prepared by the 5' addition of ³²P using [γ -³²P]ATP (3000 Ci/ μ mol; 1 Ci = 37 GBq; New England Nuclear) and T4 polynucleotide kinase (New England Nuclear) (29).

Screening of cDNA Clones. Recombinant cDNA clones were selected on LB agar plates supplemented with 12.5 μ g of tetracycline per ml. Nitrocellulose replicas of the cDNA library were prepared and replica clones were screened by *in situ* hybridization with ³²P-labeled mixed oligonucleotides (29, 31).

Plasmid DNAs prepared as described (32) from clones hybridizing with the probe oligonucleotides were further screened and insert size was determined by Southern hybrid-

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.

Abbreviation: bp, base pair(s).

ization (33) following restriction with *Pst* I (Bethesda Research Laboratories) and separation by agarose electrophoresis.

Hybridizations with ONT 17-1, ONT 17-2, and ONT 11 were carried out at 40°C, 42°C, and 22°C, respectively, as described (29).

Nucleotide Sequence Analysis. Plasmid cDNA inserts were cloned into the *Pst* I site of M13mp9 (34). Clones representing opposite cDNA orientations were selected by phage hybridization (C test) (35). Non-random BAL-31 exonuclease (New England Nuclear) deletion libraries from each orientation were prepared as described by Poncz *et al.* (36). Sequence analysis was performed by the dideoxy nucleotide chain-termination procedure of Sanger *et al.* (37) using deoxyadenosine 5'-[α -³⁵S]thio]triphosphate and a 17-nucleotide universal primer (P-L Biochemicals). DNA samples were electrophoresed on 6% acrylamide ion gradient gels (38). Data analysis was performed by using the Staden (39) DB system sequence analysis program.

RNA Transfer Blots. RNA nitrocellulose transfer blots and hybridization were done as described by Thomas (40) following electrophoresis of glyoxal-denatured poly(A)⁺ RNA on 1.2% agarose gels. Approximately 10 μ g of poly(A)⁺ RNA was applied per lane. Plasmid DNA containing the appropriate cDNA insert was labeled by nick-translation using [³²P]dCTP and a nick-translation kit (Bethesda Research Laboratories). The RNA was hybridized with nick-translated [α -³²P]DNA at 42°C in the presence of 50% formamide. Denatured *Hae* III restriction fragments of ϕ X174 (Bethesda Research Laboratories) were used as size markers. ϕ X174 *Hae* III fragments were visualized by staining that part of the gel containing the size markers with acridine orange (30 μ g/ml).

RESULTS

NH₂-Terminal Amino Acid Sequence of Rat Yolk Sac Tumor Proteoglycan and Choice of Oligonucleotides. A single NH₂-terminal amino acid sequence with a calculated initial yield of 57% based on protein assay was obtained for the yolk sac tumor proteoglycan (Fig. 1). Three oligonucleotides, two 17-mer mixed probes and one mixed 11-mer probe ending in A-G or T-C, were synthesized based on the protein sequence as shown in Fig. 1.

Selection of cDNA Clones. A cDNA library containing \approx 35,000 clones was prepared from L2 poly(A)⁺ RNA and screened by *in situ* colony hybridization and Southern blot hybridization using the mixed 17-mer oligonucleotides ONT 17-1 and ONT 17-2 shown in Fig. 1. Six clones were selected based on hybridization with either of the oligonucleotides. These clones were further analyzed by Southern hybridization of *Pst* I-restricted purified plasmid DNA to the oligonucleotides ONT 17-2 and ONT 11 (Fig. 2). Plasmid DNA from

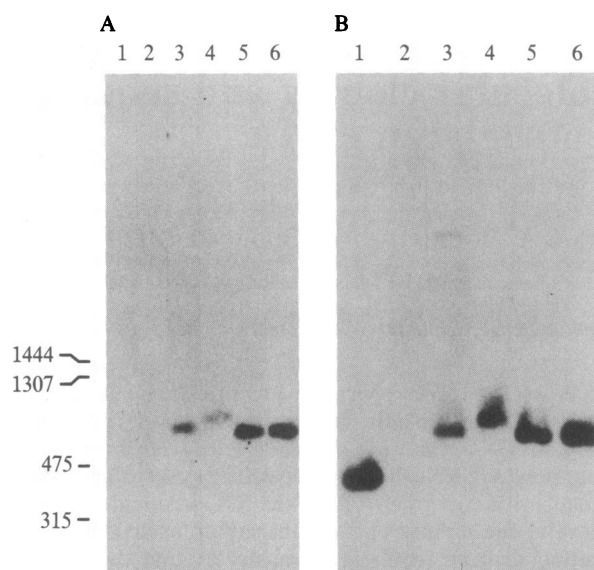


FIG. 2. Identification of yolk sac tumor proteoglycan cDNA clones and determination of insert size. Southern hybridization of plasmid DNA restricted with *Pst* I was performed by using ONT 11 (A) and ONT 17-2 (B). Lanes 1, pPG7 hybridizes with ONT 17-2 only; lanes 2, pPG8 hybridizes with ONT 17-1 only and is negative with both of the probes used here; lanes 3, pPG6; lanes 4, pPG1; lanes 5, pPG2; lanes 6, pPG3. *Taq* I fragments of pBR322 (indicated as number of bp) were used as size markers.

4 clones hybridized with both ONT 17-2 and ONT 11, whereas 2 clones hybridized with either ONT 17-1 (not shown) or ONT 17-2 alone (Fig. 2). Of the 4 clones with cDNA inserts hybridizing with both ONT 17-2 and ONT 11, the clone with the largest cDNA insert, pPG1 (lane 4), was selected for DNA sequencing.

DNA Sequence of pPG1 cDNA. DNA sequencing of both strands of pPG1 was carried out by the dideoxynucleotide chain-termination method. The 874-bp DNA sequence and inferred amino acid sequence are shown in Fig. 3. The mature proteoglycan sequence is contained within the pPG1 cDNA starting 152 nucleotides from the 5' end of pPG1 and terminating at position 463. The most distinctive feature of the cDNA coding sequence is a region coding for 25 serine residues alternating with 24 glycine residues. This region follows the 14 amino acid NH₂-terminal coding region. The remaining sequence codes for an additional 41 amino acids.

In addition to the proteoglycan coding region, pPG1 contains both 5' and 3' flanking sequences. The 5' flanking region has an open reading frame continuous with that of the proteoglycan coding region. This region does not have a recognizable signal sequence but does have two ATG codons, one immediately preceding the NH₂-terminal arginine of the

NH ₂ -terminal amino acid sequence	ONT 17-1,-2										ONT 11								
	Arg	Gly	Phe	Pro	Asn	Asp	Phe	Phe	Pro	Ile	X	Asp	Asp	Tyr	Ser	Gly	X	Gly	X
Oligonucleotide sequences																			
ONT 17-1		T	T	A	C	T	A	A	A	A	A	A	G	G	A	T	A		
ONT 17-2		T	T	A	C	T	A	A	A	A	A	A	G	G	G	C	T	A	
ONT 11		C	T	A	C	T	A	A	T	A	G	T	C						

FIG. 1. NH₂-terminal amino acid sequence of yolk sac tumor chondroitin sulfate proteoglycan and oligonucleotide probes based on the amino acid sequence are shown. Blank positions in the amino acid sequence analysis are indicated by an X. Synthetic oligonucleotides 17 and 11 nucleotides in length were derived from the two segments of the NH₂-terminal sequence indicated. Boxes indicate differences in sequence between ONT 17-1 and ONT 17-2 at position 15 and in the terminal dinucleotides of ONT 11.

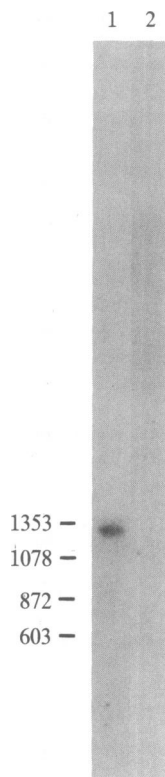


FIG. 4. RNA transfer hybridization of poly(A)⁺ RNA using ³²P-labeled nick-translated pPG1. Glyoxal-denatured rat L2 yolk sac tumor (lane 1) and rat hepatoma 7777 (lane 2) poly(A)⁺ RNA (10 μg) was electrophoresed through a 1.2% agarose gel, blotted onto nitrocellulose, and hybridized with ³²P-labeled pPG1 at 42°C in the presence of 50% formamide. Glyoxal-denatured φX174 *Hae* III fragments (indicated as number of nucleotides) were used as size markers.

glycan core proteins is accomplished via *O*-glycosyl linkage to serine (42), and all but one of the serines coded for by pPG1 are contained within the serine-glycine repeat. Moreover, it is known that glycine residues are also involved in glycosaminoglycan attachment, since glycine is abundant in proteoglycans (12, 17) and synthetic peptides containing serine and glycine can serve as acceptors for glycosaminoglycan chain initiation (43). The extent of serine *O*-glycosylation in the L2 proteoglycan has been estimated to be near 60% (17), indicating that at least 14 of the serine residues of the core protein bear a chondroitin sulfate chain.

The structure of the serine-glycine region closely parallels that predicted by Robinson *et al.* (44) for the heparin attachment region of a rat heparin proteoglycan, indicating that the serine-glycine repeat may be a general feature of at least a subset of proteoglycans. In their study, Robinson *et al.* (44) found that the Pronase-resistant heparin attachment region contained only the amino acids serine and glycine and proposed a model in which 15–20 serine residues alternate with glycine residues, with heparin chains attached to at least 2 of every 3 serines. It would appear that at least two proteoglycans, the heparin one and our rat yolk sac tumor proteoglycan, have identical or nearly identical glycosaminoglycan attachment regions. Since the codon usage for serine and glycine is quite restricted in the pPG1 cDNA, with 81% of the serine residues being coded for by two of six possible codons and 70% of the glycine residues being coded for by one of four possible codons, it is possible that the pPG1 cDNA could also hybridize to the mRNA and gene of the heparin proteoglycan and perhaps even to those of other proteoglycans.

In addition to chondroitin sulfate side chains, many pro-

teoglycans have *O*- and *N*-linked oligosaccharides (1). These oligosaccharide chains are linked to the proteoglycan through either threonine (serine) *O*-glycosyl or asparagine *N*-glycosyl acceptor recognition sequences (45). An examination of the L2 proteoglycan amino acid sequence does not reveal any asparagine oligosaccharide acceptor sites that have the sequence X-asparagine-Y-serine (45).

The mRNA sequence 5' of the NH₂ terminus of the mature proteoglycan coding sequence is in an open reading frame continuous with that of the mature proteoglycan coding region. Although there is no direct evidence that it is translated, the possibility remains that the region does code for a segment present in a proteoglycan precursor. This possibility seems somewhat unlikely, since there is a potential initiation codon immediately preceding the codon for the NH₂-terminal arginine residue of the core protein. However, if the core protein is synthesized in the form in which it appears extracellularly, it would be lacking a signal peptide sequence (46). It is unusual, but not unprecedented, for a secreted or membrane protein not to have a signal sequence (ref. 47; reviewed in ref. 48). The 5' flanking region contained in the pPG1 cDNA clone does not have a recognizable signal peptide sequence, but since the cDNA is smaller by about 425 bp than the mRNA there could be a signal sequence in the mRNA not covered by our cDNA clone. Taking into account the probable length of the poly(A) tract (150–200 bp; ref. 49) that is not represented in our cDNA, the remaining 5' sequence could be about 150–200 bp long. It is also possible that the biosynthetic pathways for proteoglycans are governed by structural rules different from those operating in other proteins and that no signal peptide is necessary for proteoglycan processing.

The identification and sequencing of the cDNA for the rat yolk sac tumor proteoglycan core protein provides a valuable tool for the further understanding of the biosynthesis of proteoglycans and the structure of their core proteins and genes. In addition, the inferred amino acid sequence provides the possibility of using synthetic peptides for generating high-titer specific proteoglycan antisera not presently available. Such antisera would be very useful in biochemical studies of proteoglycan core protein synthesis and processing.

We thank Khanh Nguyen and Yvonne Öhgrén for technical assistance. This work was supported by Grant CA 28101 and Cancer Center Support Grant CA 30199 from the National Cancer Institute.

- Hascall, V. C. & Hascall, G. K. (1981) in *Cell Biology of Extracellular Matrix*, ed. Hay, E. D. (Plenum Press), pp. 39–63.
- Kraemer, P. M. (1971) *Biochemistry* **10**, 1437–1445.
- Oldberg, A., Kjellen, L. & Hook, M. (1979) *J. Biol. Chem.* **254**, 8505–8510.
- Bumol, T. F. & Reisfeld, R. A. (1982) *Proc. Natl. Acad. Sci. USA* **79**, 1245–1249.
- Hayman, E. G., Oldberg, A., Martin, G. R. & Ruoslahti, E. (1982) *J. Cell Biol.* **94**, 28–35.
- Oohira, A., Wight, T. N., McPherson, J. & Bornstein, P. (1982) *J. Cell Biol.* **92**, 357–367.
- Perkins, M. E., Ji, T. H. & Hynes, R. O. (1979) *Cell* **16**, 941–952.
- Hayashi, M., Schlesinger, D. H., Kennedy, D. W. & Yamada, K. M. (1980) *J. Biol. Chem.* **255**, 10017–10020.
- Sakashita, S. & Ruoslahti, E. (1980) *Arch. Biochem. Biophys.* **205**, 283–290.
- Toole, B. P. (1976) *J. Biol. Chem.* **251**, 895–897.
- Laterra, J., Silbert, J. E. & Culp, L. A. (1983) *J. Cell Biol.* **96**, 112–123.
- Kanwar, Y. S., Veis, A., Kimura, J. H. & Jakubowski, M. L. (1984) *Proc. Natl. Acad. Sci. USA* **81**, 762–766.
- Brennan, M. J., Oldberg, A., Pierschbacher, M. D. & Ruoslahti, E. (1984) *J. Biol. Chem.* **259**, 13742–13750.

14. Glimelius, B., Norling, B., Westermark, B. & Wasteson, A. (1983) *Biochem. J.* **172**, 443–456.
15. Kojima, J., Nakamura, N., Kanatani, M. & Akiyama, M. (1982) *Cancer Res.* **42**, 2857–2860.
16. Robinson, J., Viti, M. & Hook, M. (1984) *J. Cell Biol.* **98**, 946–953.
17. Oldberg, A., Hayman, E. G. & Ruoslahti, E. (1981) *J. Biol. Chem.* **256**, 10847–10852.
18. Oldberg, A. & Ruoslahti, E. (1982) *J. Biol. Chem.* **257**, 4859–4863.
19. Suzuki, S., Pierschbacher, M. D., Hayman, E. G., Nguyen, K., Ohgren, Y. & Ruoslahti, E. (1984) *J. Biol. Chem.* **259**, 15307–15314.
20. Brennan, M. J., Oldberg, A., Hayman, E. G. & Ruoslahti, E. (1983) *Cancer Res.* **43**, 4302–4307.
21. Wewer, U. (1982) *Dev. Biol.* **93**, 416–421.
22. McMahon, J. B. (1978) Dissertation (University of Vermont, Burlington).
23. Ullrich, A., Shine, J., Chirgwin, J., Pictet, R., Fischer, E., Rutter, W. J. & Goodman, H. M. (1977) *Science* **196**, 1313–1319.
24. Singer, R. H. & Penman, S. (1973) *J. Mol. Biol.* **78**, 321–334.
25. Gonzalez, F. J. & Kasper, C. B. (1981) *J. Biol. Chem.* **256**, 4699–4700.
26. Suggs, S. V., Wallace, R. B., Hirose, T., Kawashima, E. H. & Itakura, K. (1981) *Proc. Natl. Acad. Sci. USA* **78**, 6613–6617.
27. Kushner, S. R. (1978) in *Genetic Engineering*, eds. Boyer, H. W. & Nicosia, S. (Elsevier/North-Holland, Amsterdam), pp. 17–23.
28. Miyoshi, K., Miyake, T., Hozumi, T. & Itakura, K. (1980) *Nucleic Acids Res.* **8**, 5461–5471.
29. Oldberg, A., Linney, E. & Ruoslahti, E. (1983) *J. Biol. Chem.* **258**, 10193–10196.
30. Tan, Z. K., Ikuta, F., Huang, T., Dugaiczky, A. & Itakura, K. (1982) *Cold Spring Harbor Symp. Quant. Biol.* **47**, 387–391.
31. Hanahan, D. & Meselson, M. (1980) *Gene* **10**, 63–67.
32. Holmes, D. S. & Quigley, M. (1981) *Ann. Biochem.* **114**, 193–197.
33. Southern, E. M. (1975) *J. Mol. Biol.* **99**, 503–512.
34. Messing, J. & Vieira, J. (1982) *Gene* **19**, 269–276.
35. Howarth, A. J., Gardner, R. C., Messing, J. & Shepherd, R. I. (1981) *Virology* **112**, 678–685.
36. Poncz, M., Solowiejczyk, D., Ballantine, M., Schwartz, E. & Surrey, S. (1982) *Proc. Natl. Acad. Sci. USA* **79**, 4298–4302.
37. Sanger, F., Nicklen, S. & Coulson, A. R. (1977) *Proc. Natl. Acad. Sci. USA* **74**, 5463–5467.
38. Biggin, M. D., Gibson, T. J. & Hong, G. F. (1983) *Proc. Natl. Acad. Sci. USA* **80**, 3963–3968.
39. Staden, R. (1982) *Nucleic Acids Res.* **10**, 27–30.
40. Thomas, P. S. (1981) *Proc. Natl. Acad. Sci. USA* **77**, 5201–5205.
41. Oldberg, A., Schwartz, C. & Ruoslahti, E. (1982) *Arch. Biochem. Biophys.* **216**, 400–406.
42. Roden, L. (1966) *J. Biol. Chem.* **241**, 5949–5954.
43. Coudron, C., Loerner, T., Jacobson, I., Roden, L. & Schwartz, N. B. (1980) *Fed. Proc. Fed. Am. Soc. Exp. Biol.* **39**, 1671 (abstr.).
44. Robinson, H. M., Horner, K. A., Hook, M., Ogren, S. & Lindahl, U. (1978) *J. Biol. Chem.* **253**, 6687–6693.
45. Hughes, R. C. (1973) *Prog. Biophys. Mol. Biol.* **26**, 189–268.
46. Blobel, G. (1975) *J. Cell Biol.* **67**, 835–851.
47. Palmiter, R. D., Gagnon, J. & Walsh, K. A. (1978) *Proc. Natl. Acad. Sci. USA* **75**, 94–98.
48. Silhavy, T. J., Benson, S. A. & Emil, S. D. (1983) *Microbiol. Rev.* **47**, 313–334.
49. Brawerman, G. (1974) *Annu. Rev. Biochem.* **42**, 621–642.