

---

**Molecular cloning and nucleotide sequence of the streptavidin gene**

---

Carlos E. Argaraña<sup>1</sup>, Irwin D. Kuntz<sup>2</sup>, Steven Birken<sup>3</sup>, Richard Axel<sup>4,5</sup> and Charles R. Cantor<sup>1</sup>

---

Departments of <sup>1</sup>Genetics and Development, <sup>4</sup>Biochemistry and Molecular Biophysics and <sup>3</sup>Medicine and <sup>5</sup>Howard Hughes Medical Research Institute, College of Physicians and Surgeons of Columbia University, 701 West 168th Street, New York, NY 10032 and <sup>2</sup>Department of Pharmaceutical Chemistry, School of Pharmacy, University of California, San Francisco, CA 94143, USA

---

Received 25 November 1985; Revised 20 January 1986; Accepted 28 January 1986

---

**ABSTRACT**

Using synthetic oligonucleotides as probes we have cloned the streptavidin gene from a genomic library of Streptomyces avidinii. Nucleotide sequence analysis indicated that a 2 Kb DNA-fragment contained the entire coding region, a signal peptide region and the 3' and 5' flanking regions of the gene. The deduced amino acid sequence shows several interrupted blocks of homology with the amino acid sequence of chicken egg-white avidin. Analysis of the secondary structure suggests a high content of beta-structure in both proteins and considerable overall structural similarity between them.

**INTRODUCTION**

Streptavidin, a 60,000 dalton protein produced by Streptomyces avidinii, forms a very strong and specific non-covalent complex with the water-soluble vitamin biotin (1-2). The protein consists of 4 identical subunits of approximate molecular weight 15,000, binds 4 mol of biotin per mol of protein, and is free of carbohydrate. Avidin, a basic glycoprotein usually isolated from chicken egg-whites, shares with streptavidin some common characteristics such as molecular weight, subunit composition and capacity to bind biotin, forming a complex of very high affinity ( $K_d \sim 10^{-15}$  M) (3,4). However the two proteins have rather different amino acid compositions, but both have an unusually high content of threonine and tryptophan.

Since biotin can be conjugated to a variety of biological molecules, the strong and specific biotin binding capacity of avidin or streptavidin has been exploited for the detection, localization or purification of proteins, carbohydrates and nucleic acids (5-8). At present, methods of biotin detection have been significantly improved and sensitive detection systems are commercially available. We were interested in cloning the

streptavidin gene to further construct expression vectors to produce fused proteins which could be easily detected or purified by means of the binding of biotin.

In this paper we describe the strategy used to clone the streptavidin gene, we present its complete nucleotide sequence, and we show the results of a comparison of the primary and secondary structure of streptavidin and avidin.

### MATERIALS AND METHODS

#### Enzymes and other reagents

All enzymes and chemicals used were from Bethesda Research Laboratories, New England Biolabs, Boehringer Mannheim Biochemicals or Pharmacia P-L Biochemicals. Radiochemicals were from New England Nuclear. Streptavidin was generously supplied by Bethesda Research Laboratories.

#### Amino acid sequence and amino acid analysis

Analysis by SDS-polyacrylamide gel electrophoresis of the preparation of streptavidin used showed, besides a main protein band, some material of lower molecular weight, possibly a degradation product of the protein. In order to obtain a pure component for amino acid sequence analysis, the preparation of streptavidin was electrophoresed in a preparative 15% slab SDS-polyacrylamide gel (9) and the main, higher molecular weight protein band, was purified from the gel. Visualization of the protein bands, elution and SDS elimination were carried out essentially according to Harger and Burgess (10).  $\text{NH}_2$ -terminal sequence analysis of the protein was performed using a Beckman 890B automatic sequencer. The identification of amino acids was carried out by HPLC (11). For amino acid analysis, the gel-purified protein was hydrolyzed with 6 N HCl in the presence of  $\beta$ -mercaptoethanol (1:1000) at 110°C under vacuo for 24 h, and the hydrolysate was analyzed on a Beckman 121MB amino acid analyzer.

#### Synthesis, purification and labelling of oligonucleotides

Oligonucleotide mixtures (see Fig. 1) were synthesized by the solid-phase phosphite triester method using an Applied Biologicals DNA/RNA synthesizer (12). The oligonucleotides were purified by preparative polyacrylamide gel electrophoresis on a

---

15% sequencing gel. Purified oligonucleotides were labelled at the 5' end with  $\gamma$ [<sup>32</sup>P]ATP (4,000-6,000 Ci/mmol) and polynucleotide kinase. Unincorporated ATP was removed by chromatography on DEAE-cellulose (13).

#### Construction of the genomic library from *Streptomyces avidinii*

Purified chromosomal DNA from *Streptomyces avidinii* was partially digested with MboI and the DNA fragments ranging between 6-19 Kb were purified by agarose gel electrophoresis. Charon 30 DNA (14) was digested to completion with BamHI, the arms isolated by agarose gel electrophoresis and then ligated with the DNA fragments of *Streptomyces avidinii* using T4 DNA ligase. The recombinant DNA was packaged in vitro into bacteriophage particles according to Maniatis et al. (15).

#### Screening of DNA clones

*E. coli* LE 392 cells were infected with the recombinant phages, plated in NZYCM-soft agarose on NZYCM (15) agar plates and grown at 37°C. Two plates containing approximately  $8 \times 10^3$  phages each were used for the screening. Three replica filters were prepared for hybridization according to Benton and Davis (16). Filters were pre-hybridized in 75 mM Tris-HCl pH 8, 100 mM sodium phosphate pH 6.5, 750 mM NaCl, 5 mM EDTA, 1% SDS, 10xDenhardt and 100 µg/ml of denatured salmon sperm DNA for 3 h at 25°C. Hybridization was done in the same solution in the presence of 4 ng/ml of labelled probe (Stvl4, see Fig 1) at a specific activity of  $10^8$ - $10^9$  cpm/µg of oligonucleotide. Filters were hybridized at 25, 28 and 31°C (one replica at each temperature) for 30-36 h then washed at 25°C for 45 min with three changes of 250 ml of the same solution used for pre-hybridization except that Denhardt and DNA were omitted. Filters were blotted dry and exposed to Kodax XR5 X-ray film with an intensifying screen.

#### DNA sequence analysis

Restriction fragments of the gene were subcloned into M13 mp18 and mp19 (17) and sequenced by the dideoxy chain termination method (18).

#### Secondary structure prediction method

Computer programs have been developed that compare the amino acid sequences of proteins to a series of sequence

patterns that have been shown to be characteristic of secondary structure elements in proteins of known tertiary structure (19-21). These patterns have been found to be approximately 90% accurate in identifying the turns that separate helices and beta strands (20). The patterns used to evaluate helical and beta propensities were taken from a study of  $\alpha/\beta$  proteins (19) augmented with others characteristic of all-helical and all-beta proteins (20). These patterns are clearly less reliable (ca. 70% correct) than the turn finding procedure. Extension of the methods to groups of sequences known to be closely related (e.g. myoglobins and immunoglobulins) did not degrade the reliability of the method (19).

### RESULTS AND DISCUSSION

#### Amino acid sequence of streptavidin

NH<sub>2</sub>-terminal amino acid analysis of a commercial preparation of streptavidin indicated the presence of both alanine and aspartic acid in the first cycle of Edman degradation of the protein. This heterogeneity can be explained by the fact that when this preparation was examined by SDS-polyacrylamide gel electrophoresis, two main protein bands with an approximate molecular weight of 17.5 and 15.5 Kd were observed. The higher molecular weight band accounted for 60-70% of the total stained protein material present in the gel. To determine the amino acid sequence, the 17.5 Kd-polypeptide chain was gel purified as described in Materials and Methods. Fig 1 shows the amino acid sequence obtained for the 40 NH<sub>2</sub>-terminal residues of the protein.

#### Isolation of the clone containing the streptavidin gene

The approach used for the isolation of the clone containing the streptavidin gene was to screen a genomic library of Streptomyces avidinii with a mixture of 16 oligonucleotides, 14 nucleotides long (Stv14) that represent all possible codon combinations for a small portion of the amino acid sequence of streptavidin (Fig 1).

Several clones, which remained positive at the three temperatures of hybridization used (see Materials and Methods) were isolated. In order to confirm the presence of the desired

Amino acid sequence determined from the gel-purified protein

```

1                               10
Asp Pro Ser Lys Asp Ser Lys Ala Gln Val Ser Ala Ala Glu Ala
                20                               30
Gly Ile Thr Gly Thr Trp Tyr Asn Gln Leu Gly Ser Thr Phe Ile
                40
Val Thr Ala Gly Ala Asp Gly Ala Leu Thr
    
```

‡ Oligonucleotide probes used

```

Amino acid sequence           7  8  9  10
Possible codons               5' AAA GCN CAA GUN 3'
Probe Stv11                   G    G
                               TTT CGN GTT CA
                               C    C

Amino acid sequence           21 22 23 24 25
Possible codons               5' UGG UAU AAU CAA CUN 3'
                               C    C    G UUU
Probe Stv14                   ACC ATA TTA GTT GA
                               G    G    C A
    
```

Figure 1. Amino-terminal amino acid sequence of streptavidin, and oligonucleotide probes used for the isolation of the streptavidin gene. (N: A,G,C and U or T).

clone, purified DNA from each presumptive positive clone was cut with BamHI, the DNA fragments were separated by agarose gel electrophoresis and analyzed by the Southern blot technique (22). In addition to Stv14, another probe, Stv11 (Fig 1) which was derived from a different part of the amino acid sequence, was used. Both probes, Stv14 and Stv11, hybridized specifically to a single fragment of approximately 2 Kb (data not shown). Nucleotide sequence analysis and amino acid sequence.

In order to identify the region containing the complementary sequence of the probe, the 2 Kb-fragment was cut with Sau3AI, subcloned into BamHI-cut M13, and the recombinants were screened with <sup>32</sup>P-labelled Stv14 probe. The DNA sequence obtained from an isolated positive clone showed the presence of part of the coding region of the gene and the sequence complementary to both probes. To localize this fragment within the 2 Kb, a partial restriction map of the 2 Kb-fragment was constructed using the method of Smith and Birnstiel (23). In order to obtain the complete nucleotide sequence of the gene, appropriate overlapping fragments were subcloned into M13 and

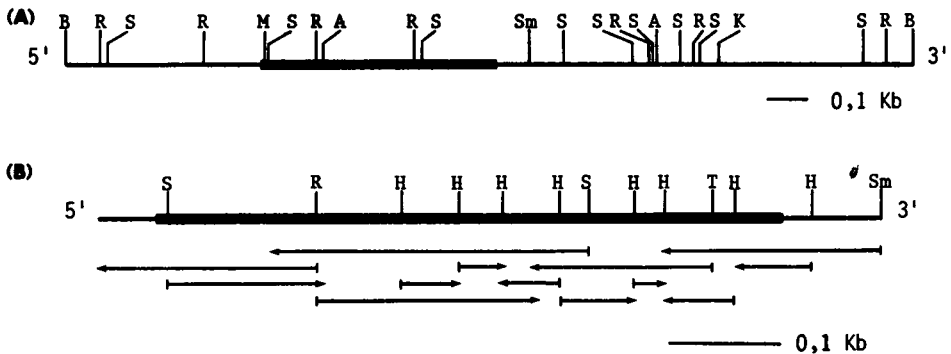


Figure 2. Partial restriction map of the cloned 2 Kb-fragment (A). Strategy used for DNA sequence analysis (B). The arrows indicate the direction and extent of the fragments sequenced. The shaded region corresponds to the coding sequence. (B: BamHI, R: RsaI, S: Sau3AI, M: MstI, A: AluI, Sm: SmaI, K: KpnI, H: HaeIII, T: Tacl).

sequenced. Fig 2 shows the partial restriction map of the 2 Kb-fragment and the strategy used to sequence the streptavidin gene.

The complete nucleotide sequence of the streptavidin gene along with the amino acid sequence is shown in Fig 3. The amino acid sequence of residues 1 to 40 is in perfect coincidence with that obtained from the protein sequence shown in Fig 1. The amino-terminal amino acid of the protein isolated in vitro is aspartic acid, thus residues -24 to -1 must be post-translationally removed to yield this mature protein. The extra 24 amino acids show common characteristics with those signal peptides present in the genes of most transmembrane and secreted proteins (24). This finding is in agreement with the fact that streptavidin has been described as a secreted protein (1). After amino-terminal processing the mature protein contains 159 amino acids and has a calculated molecular weight of 16,500 dalton, in close agreement with the value of approximately 17,500 dalton found for the streptavidin subunits by SDS-polyacrylamide gel electrophoresis (data not shown).

A comparison of the amino acid composition of streptavidin obtained from the amino acid sequence derived from the nucleotide sequence, the amino acid analysis of the gel-purified



Table 1  
Amino acid composition of streptavidin

Amino acid	Residues per subunit		
	Nucleotide <sup>a</sup> sequence	Amino acid <sup>b</sup> analysis (this work)	Amino acid <sup>c</sup> analysis (earlier work)
Lys	8	8.7	4
His	2	2.6	2
Arg	4	3.0	4
Asp	8		
Asn	10	18.0*	12*
Thr	19	18.3	19
Ser	14	13.0	10
Glu	5		
Gln	6	11.3*	9*
Pro	4	3.7	2
Gly	18	20.6	17
Ala	25	25.0	17
Cys	0	0	0
Val	10	10.1	7
Met	0	0	0
Ile	4	4.0	3
Leu	8	8.5	8
Tyr	6	6.1	6
Phe	2	2.1	2
Trp	6	4.0#	8

- (a) The composition of the mature protein after NH<sub>2</sub>-terminal processing is given.
- (b) The values were calculated from the amino acid analysis of the gel-purified protein.
- (c) The values were taken from reference (4).
- (\*) Because acid hydrolysis of proteins results in deamination of asparagine and glutamine, these amino acids are not distinguished from aspartate and glutamate.
- (#) Tryptophan recovery was low since HCl hydrolysis was employed (addition of  $\beta$ -mercaptoethanol permitted some recovery of tryptophan).

acids. It is interesting that identical or similar values are found for those amino acids that are absent or rarely present in the N- or C-terminal region of the processed streptavidin. In addition to this observation we found that a different commercial preparation of streptavidin showed a lower and variable molecular weight than the polypeptide that we used to determine the amino acid sequence. This suggests that the N- and/or C-terminal regions of the protein may be particularly susceptible to proteolytic degradation. We calculated that the 10-12 N-terminal residues plus the 19-21 C-terminal residues



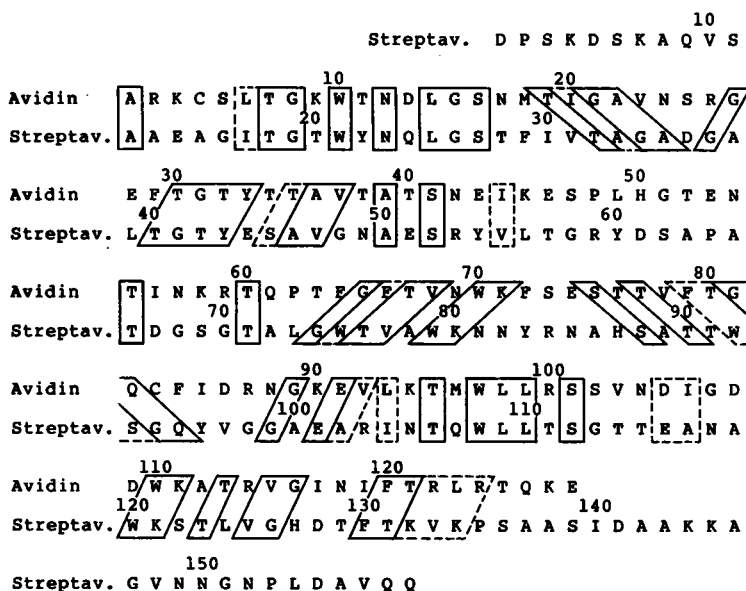


Figure 4. Amino acid sequence comparison of streptavidin and avidin. Identical residues are enclosed by solid lines and chemically similar residues by broken lines. Both sequences were aligned to give maximum homology. (Heterogeneity in residue number 34 of avidin has been reported (25); Ile or Thr is present in this position).

account, approximately, for the discrepancy found in amino acid content shown in Table 1. If this speculation is correct the previously reported amino acid analysis was probably obtained from a partially degraded streptavidin.

#### Primary and secondary structure comparison of streptavidin and avidin

Fig 4 shows the amino acid sequence of streptavidin compared with that of avidin (25), the biotin-binding protein from chicken egg-white. Streptavidin has 159 amino acids compared with 128 for avidin. Several regions of extensive homology were found between both proteins. Of particular interest is the homology around and including tryptophans 21, 79 and 120 of streptavidin. In avidin, the corresponding tryptophans 10, 70 and 110 are protected by biotin from oxidizing agents, suggesting that these residues are implicated in the biotin-binding site of the protein (4). Besides this, a

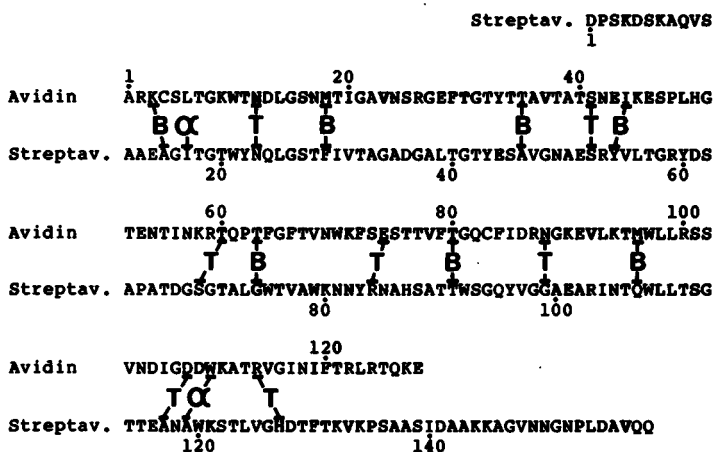


Figure 5. Comparison of predicted secondary structures of streptavidin and avidin. The sequences have been aligned as in Fig 4. α: alpha-helix, B: beta-strand, T: turn. (The final 20 C-terminal amino acids of streptavidin were not analyzed).

unique NH<sub>2</sub>-group, probably one of the three lysine residues (9, 71 and 111) which are adjacent to the tryptophans, has been found to be important for the biotin-binding activity of avidin (4). In streptavidin, two of these three lysines are conserved as lysine residues (80 and 121) also next to tryptophans.

Secondary structures were calculated for both proteins using algorithms that predict conformation from amino acid sequence (19-21). Fig 5 shows the residues at which alpha helical, beta-strand or turn features are centered. Both proteins show a clear structural homology with a high preponderance of beta structure. The alternating hydrophobic, hydrophilic pattern for most of the suggested beta-strands is consistent with a folded beta-sheet or beta-barrel geometry (26). The overall composition pattern of both sequences suggests that both proteins fall in the family of "all beta" proteins (27). The list of turns shown in Fig 5 is incomplete but there is a good probability (19) that the assigned ones are correct. The extent and exact location of beta-structure is more difficult to predict. On the other hand it is clear there is little, if any, alpha-helix in both proteins. The best chance for finding alpha-helices is in the N-terminal region of

streptavidin and the C-terminal region in both proteins.

In agreement with these predictions, avidin has been found to have a content of 55% beta-structure and 5% alpha-helix as determined by raman spectroscopy (28).

Even though both proteins show some physico-chemical differences such as isoelectric point, content of carbohydrates and amino acid composition, they have similar subunit composition, molecular weight and affinity for biotin. This, along with the homology found in several stretches of the amino acid sequence and the overall secondary structure similarity of both proteins suggests that functional and structural constraints have been remarkably conserved during their evolution. It is reasonable to speculate that there is only one way to create a binding site with such high affinity for biotin, and when the tertiary structures of both proteins are available the binding sites residues will be identical.

#### ACKNOWLEDGMENTS

We wish to thank Dan R. Littman and Jim Roberts for help with the construction of the Streptomyces avidinii library, Carlos Barreda for technical assistance, Frank Morgan for discussion and advice on the amino acid sequence of streptavidin and Wilma Saffran for her help with the oligonucleotide purification and nucleotide sequencing. C.E. Argaraña is a recipient of a postdoctoral fellowship from the Consejo Nacional de Investigaciones Cientificas y Tecnicas of Argentina. This work was supported, in part, by grants from the National Institute of Health, GM14825, HD15454 and the National Cancer Institute, CA39782.

#### REFERENCES

1. Chalet, L., Miller, T.W., Tausing, F. and Wolf, F.J. (1963) *Antimicrob. Agents Chemother.* 3, 28-32.
2. Chalet, L. and Wolf, F.J. (1964) *Arch. Biochem. Biophys.* 106, 1-5.
3. Green, N.M. (1963) *Biochem. J.* 89, 585-589.
4. Green, N.M. (1975) *Advances in Protein Chemistry* 29, 85-133.
5. Bayer, E.A. and Wilcheck, M. (1980) *Methods of Biochemical Analysis* 26, 1-46.
6. Langer, R.L., Waldrop, A.A. and Ward, D.C. (1981) *Proc. Natl. Acad. Sci. USA* 78, 6633-6637.

7. Haeuptle, M.-T., Aubert, M.L., Djiane, J.-P. and Kraehenbuhl, J. (1983) *J. Biol. Chem.* 258, 305-314.
8. Hsu, S.-M. and Raine, L. (1981) *J. Histochem. Cytochem.* 29, 1349-1353.
9. Laemmli, U. (1970) *Nature* 227, 680-685.
10. Hager, D.A. and Burgess, R.R. (1980) *Anal. Biochem.* 109, 76-86.
11. Zimmerman, C.L., Apella, E. and Pisano J.J. (1973) *Anal. Biochem.* 77, 569-573.
12. Urbina, G.A., Sathe, G.M., Liu, W., Guillen, M.F., Duck, P.D., Bender, R. and Ogilvie, K.K. (1981) *Science* 214, 270-274.
13. Wallace, R.B. (1981) *Gene* 16, 21-26.
14. Rimm, D.V., Horness, D., Kucera, J. and Blattner, F.R. (1980) *Gene* 12, 301-309.
15. Maniatis, T., Fritsch, E.F. and Sambrook, J. (1982) *Molecular Cloning.* ( Cold Spring Harbor, New York: Cold Spring Harbor Laboratory).
16. Benton, W.D. and Davis, R.W. (1977) *Science* 196, 180-182.
17. Vieira, J. and Messing, J. (1982) *Gene* 19, 259-268.
18. Sanger, F., Nicklen, S. and Coulson, A.R. (1977) *Proc. Natl. Acad. Sci. USA* 74, 5463-5467.
19. Cohen, F.E., Abarbanel, R.M., Kuntz, I.D. and Fletterick, R.J. (1983) *Biochemistry* 22, 4894-4904.
20. Cohen, F.E., Abarbanel, R.M., Kuntz, I.D. and Fletterick, R.J. (1986) *Biochemistry* (in press).
21. Abarbanel, R.M. (1984) PhD thesis, University of California.
22. Southern, E.M. (1975) *J. Mol. Biol.* 98, 503-517.
23. Smith, H.O. and Birnstiel, M.L. (1976) *Nucleic Acids Res.* 3, 2387-2398.
24. Kreil, G. (1981) *Ann. Rev. Biochem.* 50, 317-348.
25. De Lange, R.J. and Huang, T.-S. (1971) *J. Biol. Chem.* 246, 698-709.
26. Richardson, J.S. (1981) *Adv. Protein Chem.* 34, 167-339.
27. Sheridan, R. Dixon, J.S., Venkataraghavan, R., Scott, K. and Kuntz, I.D. (1985) *Biopolymers* 24, 1995-2003.
28. Honzatko, R.B. and Williams, R.W. (1982) *Biochemistry* 21, 6201-6205.