# Amino acid sequence of human plasma $\alpha_1$B-glycoprotein: Homology to the immunoglobulin supergene family

(protein structure/gene duplication/secretory component/immunoglobulin receptor)

NORIAKI ISHIOKA, NOBUHIRO TAKAHASHI, AND FRANK W. PUTNAM*

Department of Biology, Indiana University, Bloomington, IN 47405

*Contributed by Frank W. Putnam, December 13, 1985*

ABSTRACT    The complete amino acid sequence has been determined for $\alpha_1$B-glycoprotein ($\alpha_1$B), a protein of unknown function present in human plasma. This protein ($M_r \approx 63,000$) consists of a single polypeptide chain N-linked to four glucosamine oligosaccharides. The polypeptide has five intrachain disulfide bonds and contains 474 amino acid residues. Analysis of the amino acid sequence by several computer programs shows that $\alpha_1$B exhibits internal duplication and consists of five repeating structural domains, each containing about 95 amino acids and one disulfide bond. $\alpha_1$B has a unique amino acid sequence. However, several domains of $\alpha_1$B, especially the third, show statistically significant homology to variable regions of certain immunoglobulin light and heavy chains. $\alpha_1$B also exhibits sequence similarity to other members of the immunoglobulin supergene family such as the receptor for transepithelial transport of IgA and IgM and the secretory component of human IgA. Because of its internal duplication and its sequence homology to immunoglobulin-like proteins, $\alpha_1$B appears to have evolved from an ancestral gene similar to that of the immunoglobulin supergene family.

We present the complete amino acid sequence of human $\alpha_1$B-glycoprotein ($\alpha_1$B) and show that it is homologous in structure to certain domains in immunoglobulins and in the receptor for polymeric immunoglobulins (poly-IgR). $\alpha_1$B was described by Schultze *et al.* (1) as an "easily precipitable $\alpha_1$-glycoprotein" present in human plasma and was later shown to be the same as the $\alpha_1$B-glycoprotein observed by Burtin (2) on immunoelectrophoresis of serum. $\alpha_1$B has been reported to have a molecular weight of 68,000 and a carbohydrate content of 13.3% (3). The polypeptide structure was puzzling because $\alpha_1$B appeared to exist in serum in two molecular forms, one having a single polypeptide chain ($M_r = 68,000$), the other seeming to have two subunits ($M_r \approx 50,000$ and 20,000). Like most plasma glycoproteins, $\alpha_1$B exhibits electrophoretic heterogeneity near its isoelectric point (pH 4.4–4.6). $\alpha_1$B is present in normal adult serum at an average concentration of 22 mg/dl; however, no change in disease has been observed nor has any biological function been proposed (3). No information on the amino acid sequence of $\alpha_1$B has been reported previously. Thus, $\alpha_1$B-glycoprotein is one of a series of human plasma glycoproteins of unidentified physiological function that have been highly purified and have been characterized by physicochemical methods but whose primary structure was unknown (3, 4).

We are engaged in a program of study of such proteins with the objectives of determining their primary structures, their relationships to other plasma proteins, and their possible functions. Earlier we reported the complete amino acid sequence of $\beta_2$-glycoprotein I (5), ceruloplasmin (6), hemopexin (7), and leucine-rich $\alpha_2$-glycoprotein (8). A nota-

ble structural feature of these four proteins, which is shared by many other plasma proteins (9, 10), is a pattern of internal duplication in amino acid sequence that is individually characteristic and is highly significant statistically. Furthermore, by computer analysis of their sequences many plasma proteins can be grouped into families that are homologous in structure, and unexpected relationships sometimes are revealed (9). For example, blood coagulation factors VIII (11) and V (12), which are deficient or defective in certain hemophilias, were recently shown to exhibit internal duplication and to have a surprising homology in amino acid sequence to ceruloplasmin, a copper oxidase. Likewise, $\beta_2$-glycoprotein I (5) is unexpectedly homologous to several proteins of the alternative complement pathway, such as the C4b binding protein and also factor B, which is a serine protease (13). Even more surprising is the homology between $\beta_2$-glycoprotein I and the human interleukin-2 receptor (14) that we have identified in separate studies.

In this investigation we determined the complete amino acid sequence of human $\alpha_1$B and found that it consists of a single polypeptide chain of 474 amino acids with four glucosamine oligosaccharides. Computer analysis of the sequence showed that $\alpha_1$B exhibits internal duplication and consists of five repeating structural domains, each containing 92–98 residues. Some of the domains of $\alpha_1$B show significant homology to variable (V) and constant (C) regions of certain immunoglobulins. Likewise, there is statistically significant homology between $\alpha_1$B and the secretory component (SC) of human IgA (15) and also with the extracellular portion of the rabbit receptor for transepithelial transport of polymeric immunoglobulins (IgA and IgM). Mostov *et al.* (16) have called the latter protein the poly-Ig receptor or poly-IgR and have shown that it is the precursor of SC. These results suggest that $\alpha_1$B belongs to the immunoglobulin supergene family—i.e., the group of proteins that have immunoglobulin-like domains, including histocompatibility antigens, the T-cell antigen receptor, poly-IgR, and other proteins involved in the vertebrate immune response (17).

## MATERIALS AND METHODS

**Materials.** Purified $\alpha_1$B prepared from human serum and antiserum to $\alpha_1$B were obtained from Behringwerke. The protein was judged to be pure by NaDodSO$_4$/polyacrylamide gel electrophoresis in both the presence and the absence of 2-mercaptoethanol, by immunodiffusion, and by automated sequence analysis of the intact $\alpha_1$B.

**Methods.** The primary structure of $\alpha_1$B was determined by methods described previously (5–8). The purified $\alpha_1$B was reduced and carboxymethylated before sequence analysis. The carboxymethylated protein was subjected to separate digestions with L-1-tosylamido-2-phenylethyl chloromethyl

Abbreviations: $\alpha_1$B, $\alpha_1$B-glycoprotein; poly-IgR, receptor for polymeric immunoglobulins; SC, secretory component.
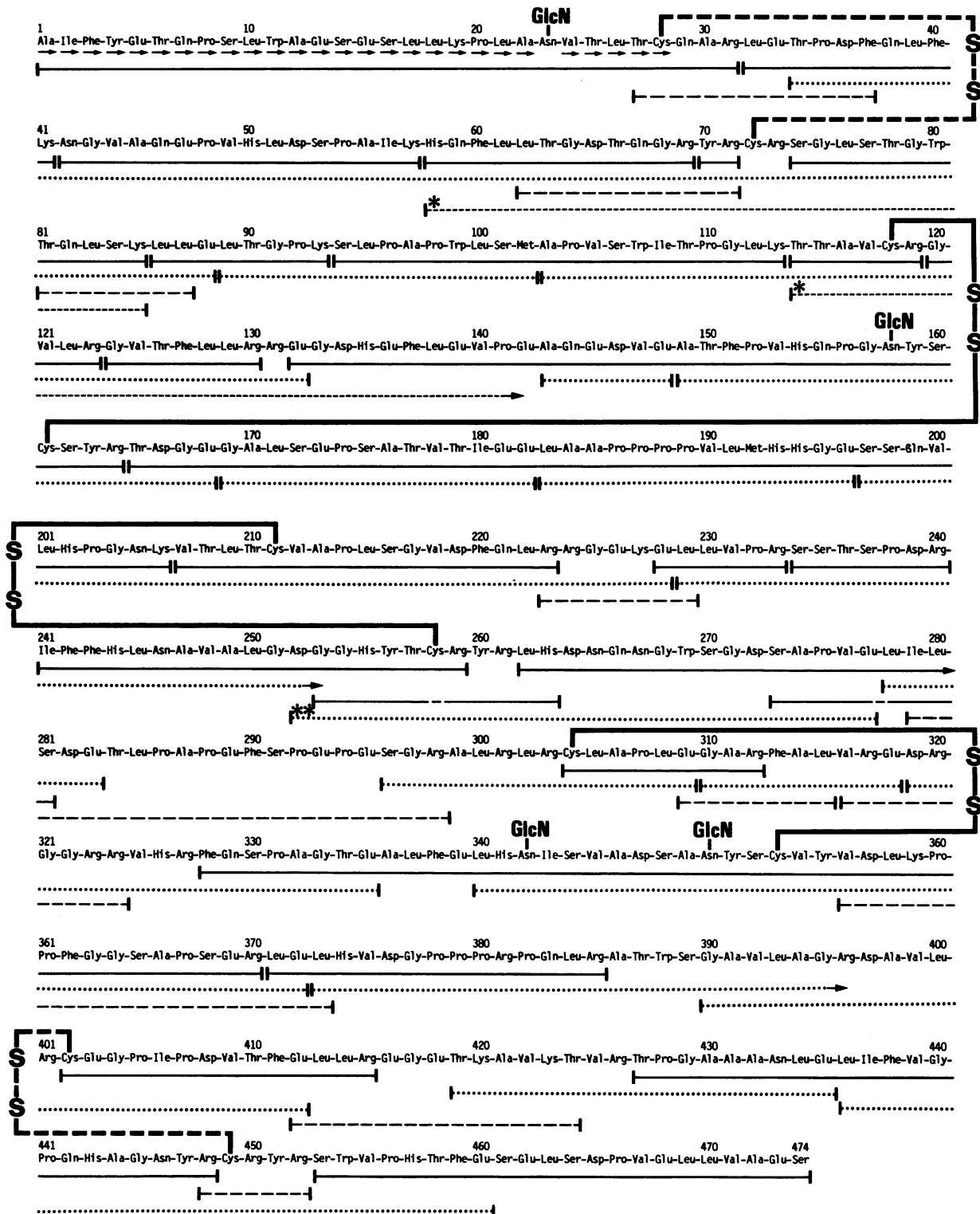*To whom reprint requests should be addressed.

FIG. 1.    Summary of the complete amino acid sequence of human $\alpha_1$B. The sequence is shown along with all peptides necessary for the proof of sequence; however, many other peptides were sequenced, and all but a few minor peptides support this structure. The peptides obtained from different digestions are as follows: —, tryptic peptides; ···, *S. aureus* V8 peptides; ---, chymotryptic peptides. A single asterisk denotes endoproteinase Lys-C peptides of a CNBr fragment. The double asterisk indicates a V8 peptide of a CNBr fragment. The peptides from Gly-253 through His-263 and from Ser-273 through Ser-281 are from a dilute acid digest. The arrows at the amino terminus indicate sequenator analysis of the intact $\alpha_1$B. GlcN indicates a glucosamine attachment site. Disulfide bonds that were proven are shown by solid lines; those that are probable are indicated by broken lines.

ketone-treated trypsin, *Staphylococcus aureus* V8 protease, and chymotrypsin. Three CNBr fragments (amino-terminal portion, middle portion, and carboxyl-terminal portion) of $\alpha_1$B that were separated by gel filtration were also digested

by endoproteinase Lys-C and also cleaved with dilute 0.07 M HCl at 108°C for 12 hr. Each enzymatic digest was separated by a combination of gel filtration and high-performance liquid chromatography (HPLC) (18). The digest with chymotrypsin

Table 1.  Amino acid composition of human α1B based on the complete amino acid sequence determination

| Amino acid | No. of residues | Amino acid | No. of residues |
|---|---|---|---|
| Aspartic acid | 19 | Valine | 34 |
| Asparagine | 11 | Methionine | 2 |
| Threonine | 29 | Isoleucine | 9 |
| Serine | 36 | Leucine | 57 |
| Glutamic acid | 39 | Tyrosine | 10 |
| Glutamine | 15 | Phenylalanine | 18 |
| Proline | 42 | Lysine | 11 |
| Glycine | 37 | Histidine | 15 |
| Alanine | 39 | Arginine | 34 |
| Half-cystine | 10 | Tryptophan | 7 |

The molecular weight of the unmodified polypeptide chain is 51,940; the number of residues is 474. Asn-23, Asn-158, Asn-342, and Asn-350 are linked to glucosamine oligosaccharides.

was separated by an automated tandem HPLC system (19). The digest of the CNBr middle portion of α1B with endoproteinase Lys-C was separated by ion-exchange HPLC (SynChropak AX-300 column, SynChrom, Linden, IN). The purified peptides were analyzed with the Beckman model 121M amino acid analyzer, and their sequences were determined by automatic Edman degradation with the Beckman model 890C sequencer (5–8). Hexosamine analysis was also done with the amino acid analyzer, after acid hydrolysis.

**Computer Analysis of Sequence Data.** The Protein Sequence Database of the Protein Identification Resource (formerly the *Atlas of Protein Sequence and Structure*), which was updated to August 1985, and the programs SEARCH, ALIGN, RELATE, PRPLOT, and DOTMATRIX were provided by the National Biomedical Research Foundation[†]. The programs SEARCH, RELATE, and ALIGN were used either with the unitary matrix or with the mutation data matrix; all gave a score for statistical significance in standard deviations (SD) of the real score above a score of 100 random runs.

## RESULTS AND DISCUSSION

**Amino Acid Composition, Polypeptide Structure, and Molecular Weight.** Human α1B consists of a single polypeptide

---

chain containing 474 amino acid residues and four glucosamine oligosaccharides (Fig. 1). Except for the high content of leucine (12.0 mol %) there is nothing notable about the amino acid composition (Table 1). It is difficult to determine the $M_r$ of glycoproteins accurately. The $M_r$ calculated from the amino acid content of the polypeptide chain is 51,940. The four glucosamines would add 10,000 to 12,000 to give a calculated $M_r$ of about 63,000 ± 1000, which is close to the reported value of 68,000 (3). In the Weber and Osborn method of NaDodSO4 electrophoresis (20), we estimated the $M_r$ as 67,000, and α1B migrated a little more slowly than bovine serum albumin. However, NaDodSO4 electrophoresis tends to give $M_r$ values for glycoproteins that are too high. In the Laemmli method (21) the estimated $M_r$ of the intact glycoprotein was 80,000, and it decreased to 68,500 after treatment with N-glycanase to remove the carbohydrate. In view of this we conclude that the $M_r$ of α1B is about 63,000 ± 1000.

**Number and Location of Oligosaccharides.** Fig. 1 gives the complete amino acid sequence of α1B and shows the linkage sites of the carbohydrate. All four glucosamine oligosaccharides are attached to asparagine in the obligate tripeptide acceptor sequence Asn-Xaa-Ser/Thr, in which Xaa is almost any amino acid. Although there is no pattern to the location of the sites, two of the oligosaccharides are linked to homologous asparagines in the duplicated tetrapeptide sequence Asn-Tyr-Ser-Cys (Asn-158 and Asn-350 in Fig. 1). The cysteine participates in a disulfide bond, but the carbohydrate must be on the surface of the molecule. Although nothing has been reported about the structure of the oligosaccharides of α1B, the most common carbohydrate in human plasma glycoproteins is a complex dibranched glucosamine oligosaccharide (22).

**Disulfide Bridges.** α1B contains 10 cysteine residues, which appear to be linked in five homologous disulfide bridges. The location of the disulfide bonds was established by tryptic digestion of the unreduced protein and purification of the disulfide-linked peptides. The latter were first separated by gel filtration in which the bridged peptides elute at a different position than the unlinked peptides from the reduced and alkylated protein. After further purification by HPLC, each bridged peptide was sequenced directly without reduction. The linkages were clearly established for the second, third, and fourth disulfide bonds, which are shown as solid lines in Fig. 1. The first cysteine (Cys-28) was found to be linked to the tryptic dipeptide Cys-Arg. However, this dipeptide could also have been derived from the similar sequence around
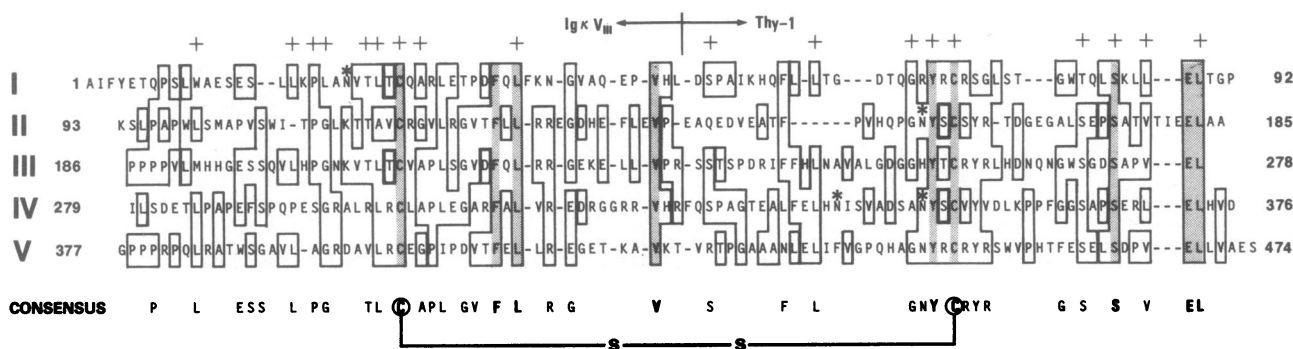


FIG. 2.  Internal homology in the primary structure of α1B. The sequence of Fig. 1 with position numbers is shown in the one-letter code for amino acids (23) and is aligned into five homologous segments or domains (I to V). Amino acids that occur five times in a vertical column are enclosed in shaded boxes in the figure; these are shown in bold type in the consensus sequence, which also contains others that occur three or four times. The homologous intradomain disulfide bond is shown in the consensus sequence. The plus signs above the figure on the left denote invariant residues in the amino-terminal half of the sequences for immunoglobulin κ light chains of subgroup III (IgκVIII). Plus signs on the right designate residues in the carboxyl-terminal half of the sequence of the mouse neuronal cell Thy-1 antigen (24) that match the consensus sequence of α1B.
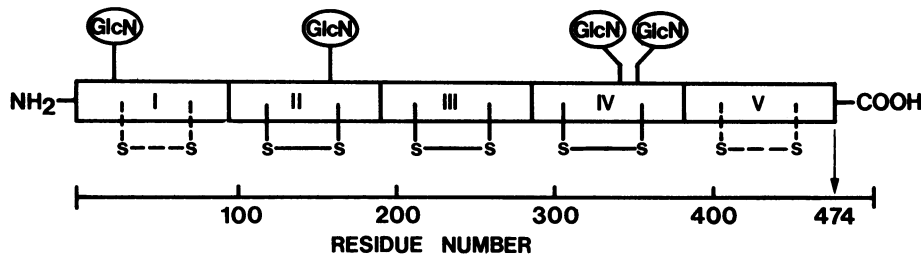
FIG. 3. Linear model of the primary structure of $\alpha_1$B. The locations of the glucosamine oligosaccharides and of the disulfide bonds are based on Fig. 1 and shown to the scale below. The division of the domains, numbered I to V, is based on Figs. 1 and 2.

Cys-449. Hence, the first and the fifth disulfide bridges are less certain and are shown with broken lines in Fig. 1.

**Internal Homology of $\alpha_1$B.** Like many plasma glycoproteins, $\alpha_1$B exhibits significant internal homology in amino acid sequence that is indicative of a series of similar structural domains. The computer program RELATE aligned the entire sequence of $\alpha_1$B into five homologous domains ranging in length from 92 to 98 residues, with an average length of 95 (Fig. 2). The intersegment homology is highly significant statistically. In SD units the scores for the paired domains range from 3.84 (domains I and II) to 8.83 (domains IV and V). Most pairings score 6 to 7 SD units, and a score of 3.0 is considered significant (23).

The homology of the domains is greatest around the cysteine residues; this has structural import because of the homologous disulfide bond formed between the pair of cysteines in each domain. Furthermore, there is one pair of identical pentapeptide sequences that involve the first cysteine in the domain (see Val-24 and Val-207 in Fig. 2), and a pair of identical tetrapeptide sequences includes the second cysteine (see Asn-158 and Asn-350). Also, 9 of the 10 tyrosines in the molecule are clustered around the carboxyl-terminal cysteines of the intradomain disulfide bridges. Altogether, in the alignment of Fig. 2 nine positions have identical residues in all five domains, and about half the residues are matched at homologous positions in two or more domains. Thus, $\alpha_1$B has about the same degree of internal homology in sequence as do the constant regions of immunoglobulin heavy chains. The consensus sequence of $\alpha_1$B based on three or more identical residues in homologous position is given in Fig. 2. In the evolutionary development of the $\alpha_1$B gene, domain III may be the closest to the primordial building block of about 95 residues because its statistical score when compared to the other four domains is consistently about 7 SD units.

**Model of the $\alpha_1$B Structural Domains.** The internal homol-

ogy in primary structure described above and the presence of an intrasegment disulfide bond suggest that $\alpha_1$B is composed of five structural domains that arose by duplication of a primordial gene coding for about 95 amino acid residues. Fig. 3 diagrams the domain structure of $\alpha_1$B based on this conclusion. Although the primary structure and the disulfide bridges follow a repeating pattern, only two of the glucosamine oligosaccharides are in homologous positions (see Asn-158 and Asn-350 in Fig. 2). However, in immunoglobulins and other plasma glycoproteins composed of repeating domains, the glucosamine carbohydrate also is not located at internally homologous positions—probably because its presence is dictated by a signal sequence. Unlike immunoglobulins (25), ceruloplasmin (6), and hemopexin (7), $\alpha_1$B is not subject to limited interdomain cleavage by proteolytic enzymes. At least, we were not able to produce such fragments by use of a variety of proteases. This stability of $\alpha_1$B is probably associated with the frequency of proline in the sequences linking the domains (Fig. 2).

**Sequence Homology to Other Proteins.** To examine the possible occurrence in other proteins of sequences homologous to $\alpha_1$B, we used the SEARCH program to compare 30-residue segments of $\alpha_1$B to the entire PIR-NBRF database. $\alpha_1$B was found to have a unique amino acid sequence. However, certain segments of immunoglobulins and related proteins gave a significant score for homology. When the ALIGN program was used to compare entire domains of $\alpha_1$B with a series of domains of immunoglobulins and related proteins, several domains of $\alpha_1$B, especially domain III, appeared to exhibit significant homology to certain members of the immunoglobulin supergene family. Scores of about 3.0 SD units or higher were given by rabbit poly-IgR, human SC, certain randomly selected V regions of human $\kappa$ and $\lambda$ light chains (V$_L$) and of human heavy chains (V$_H$), and a few domains of the constant regions of the five
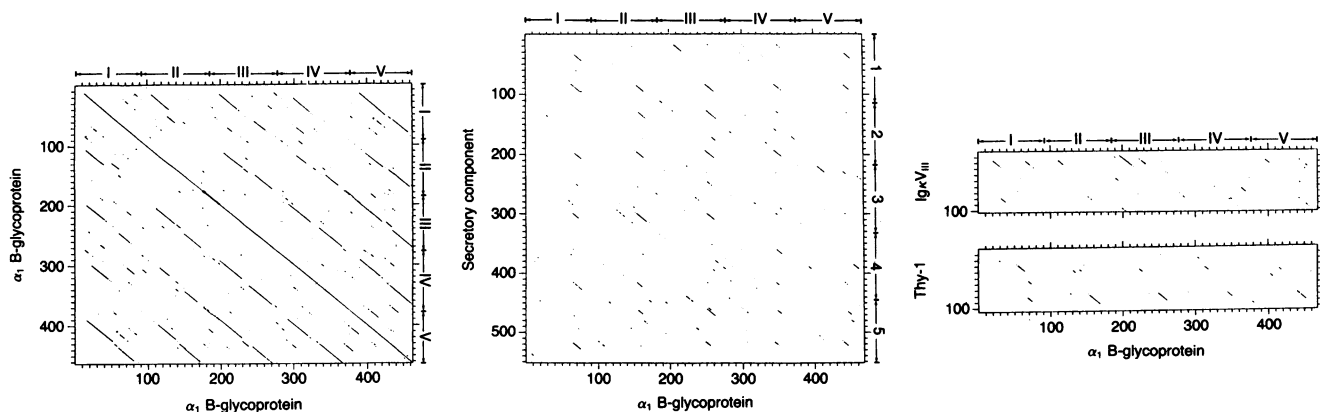


FIG. 4. Graphic matrix plots generated by the DOTMATRIX computer program of comparisons of the amino acid sequence of $\alpha_1$B and proteins of the immunoglobulin supergene family. $\alpha_1$B is compared to itself (*Left*), to human SC (*Center*), to the V region of human $\kappa$ light chains of subgroup III (*Upper Right*) and the mouse neuronal cell Thy-1 antigen (*Lower Right*). The five domains of $\alpha_1$B are numbered I to V. For $\alpha_1$B the slanting lines paralleling the solid diagonal line indicate extensive sequence similarity of the domains of $\alpha_1$B to each other, signifying a 5-fold internal repetition. For the other proteins the sets of short slanting lines reflect the fact that sequence similarity is greatest around the disulfide bridges.

human heavy chain classes ($C_H$). Domain III of $\alpha_1B$ frequently scored 3–6 SD units when compared to the $V_L$ and $V_H$ regions even though the immunoglobulin domain size is about 110 amino acid residues compared to 95 in $\alpha_1B$. Domain III generally scored higher than other $\alpha_1B$ domains when compared with immunoglobulin $C_H$ domains, but for most $C_H$ domains the sequence homology with $\alpha_1B$ was marginal or not significant. One-third of the comparisons between $\alpha_1B$ domains and poly-IgR domains gave scores of about 3–4 SD units, with domain III generally scoring higher than the others. Surprisingly, $\alpha_1B$ appeared less homologous to human SC than to rabbit poly-IgR. The score for comparison of the entire 474 residues of $\alpha_1B$ to the extracellular portion of poly-IgR (residues 30–558) was significant at the level of 4.8 SD units but was at the marginal level of 2.3 SD units when the whole $\alpha_1B$ protein was compared to the whole SC protein (residues 1–558).

The presence of five homologous domains in $\alpha_1B$ and their relationship to immunoglobulin-like domains was confirmed by sequence comparisons generated as graphic matrix plots by the DOTMATRIX program (Fig. 4). These results and those described above suggest that $\alpha_1B$ has immunoglobulin-like domains, including the intradomain disulfide bridge, and may be a member of the immunoglobulin supergene family. Because of its homology to membrane-bound poly-IgR and to SC (the plasma protein derived by proteolytic cleavage of poly-IgR), $\alpha_1B$ may be related to the long-sought membrane receptor for IgG. Not only is $\alpha_1B$ homologous to the immunoglobulin supergene family in amino acid sequence, but also, like immunoglobulins and their receptors, it appears to have evolved by gene duplication from a primordial gene coding for a structural domain of similar size.

1. Schultze, H. E., Heide, H. & Haupt, H. (1963) *Nature (London)* **200**, 1103.
2. Burtin, P. (1964) in *Immunoelectrophoretic Analysis*, eds. Grabar, P. & Burtin, P. (Elsevier, Amsterdam), pp. 94–124.
3. Schwick, H. G. & Haupt, H. (1984) in *The Plasma Proteins*, ed. Putnam, F. W. (Academic, Orlando, FL), 2nd Ed., Vol. 4, pp. 167–220.
4. Putnam, F. W. (1984) in *The Plasma Proteins*, ed. Putnam, F. W. (Academic, Orlando, FL), 2nd Ed., Vol. 4, pp. 45–166.
5. Lozier, J., Takahashi, N. & Putnam, F. W. (1984) *Proc. Natl. Acad. Sci. USA* **81**, 3640–3644.
6. Takahashi, N., Ortel, T. L. & Putnam, F. W. (1984) *Proc. Natl. Acad. Sci. USA* **81**, 390–394.
7. Takahashi, N., Takahashi, Y. & Putnam, F. W. (1985) *Proc. Natl. Acad. Sci. USA* **82**, 73–77.
8. Takahashi, N., Takahashi, Y. & Putnam, F. W. (1985) *Proc. Natl. Acad. Sci. USA* **82**, 1906–1910.
9. Doolittle, R. F. (1984) in *The Plasma Proteins*, ed. Putnam, F. W. (Academic, Orlando, FL), 2nd Ed., Vol. 4, pp. 317–359.
10. Putnam, F. W. (1985) *Protides Biol. Fluids Proc. Colloq.* **25**, 407–410.
11. Vehar, G. A., Keyt, B., Eaton, D., Rodriguez, H., O'Brien, D. P., Rotblat, F., Opperman, H., Keck, R., Wood, W. I., Harkins, R. N., Tuddenham, E. G. D., Lawn, R. W. & Capon, D. J. (1984) *Nature (London)* **312**, 337–342.
12. Church, W. R., Jernigan, R. L., Toole, J., Hewick, R. M., Knopf, J., Knutson, G. J., Nesheim, M. E., Mann, K. G. & Fass, D. N. (1984) *Proc. Natl. Acad. Sci. USA* **81**, 6934–6937.
13. Chung, L. P., Bentley, D. R. & Reed, K. B. M (1985) *Biochem. J.* **230**, 133–141.
14. Leonard, W. J., Depper, J. M., Kanehisa, M., Kronke, M., Peffer, N. J., Svetlik, P. B., Sullivan, M. & Greene, W. C. (1985) *Science* **230**, 633–639.
15. Eiffert, H., Quentin, E., Decker, J., Hillemeir, S., Hufschmidt, M., Klingmuller, D., Weber, M. H. & Hilschmann, N. (1984) *Hoppe-Seyler's Z. Physiol. Chem.* **365**, 1489–1495.
16. Mostov, K. E., Friedlander, M. & Blobel, G. (1984) *Nature (London)* **308**, 37–43.
17. Hood, L., Kronenberg, M. & Hunkapiller, T. (1985) *Cell* **40**, 225–229.
18. Takahashi, N., Takahashi, Y., Ortel, T. L., Lozier, J., Ishioka, N. & Putnam, F. W. (1984) *J. Chromatogr.* **317**, 11–26.
19. Takahashi, N., Ishioka, N., Takahashi, Y. & Putnam, F. W. (1985) *J. Chromatogr.* **326**, 407–418.
20. Weber, K. & Osborn, M. (1969) *J. Biol. Chem.* **244**, 4406–4412.
21. Laemmli, U. K. (1970) *Nature (London)* **227**, 680–685.
22. Baenziger, J. U. (1984) in *The Plasma Proteins*, ed. Putnam, F. W. (Academic, Orlando, FL), 2nd Ed., Vol. 4, pp. 271–315.
23. Barker, W. C., Ketcham, L. K. & Dayhoff, M. O. (1978) in *Atlas of Protein Sequence and Structure*, ed. Dayhoff, M. O. (Natl. Biomed. Res. Found., Washington, DC), Vol. 5, Suppl. 3, pp. 359–362.
24. Williams, A. F. & Gagnon, J. (1982) *Science* **216**, 696–703.
25. Putnam, F. W. (1977) in *The Plasma Proteins*, ed. Putnam, F. W. (Academic, New York), 2nd Ed., Vol. 3, pp. 1–153.