# The complete amino acid sequence of human serum transferrin

(primary structure/protein evolution/contiguous gene duplication/disulfide bonds/iron transport)

Ross T. A. MacGillivray*, Enrique Mendez†, Sudhir K. Sinha, Michael R. Sutton‡, Janet Lineback-Zins, and Keith Brew§

Department of Biochemistry, University of Miami School of Medicine, Miami, Florida 33101

ABSTRACT    The complete amino acid sequence of human serum transferrin has been determined by aligning the structures of the 10 CNBr fragments. The order of these fragments in the polypeptide chain is deduced from the structures of peptides overlapping methionine residues and other evidence. Human transferrin contains 678 amino acid residues and—including the two asparagine-linked glycans—has an overall molecular weight of 79,550. The polypeptide chain contains two homologous domains consisting of residues 1–336 and 337–678, in which 40% of the residues are identical when aligned by inserting gaps at appropriate positions. Disulfide bond arrangements indicate that there are seven residues between the last half-cystine in the first domain and the first half-cystine in the second domain and therefore, a maximum of seven residues in the region of polypeptide between the two domains. Transferrin—which contains two Fe-binding sites—has clearly evolved by the contiguous duplication of the structural gene for an ancestral protein that had a single Fe-binding site and contained ≈340 amino acid residues. The two domains show some interesting differences including the presence of both N-linked glycan moieties in the COOH-terminal domain at positions 413 and 610 and the presence of more disulfide bonds in the COOH-terminal domain (11 compared to 8). The locations of residues that may function in Fe-binding are discussed.

The iron-transport protein of serum, transferrin, is a monomeric glycoprotein with $M_r \approx 80,000$. The properties and functions of serum transferrins have been recently reviewed in detail (1). Briefly, the transferrin molecule possesses two independent metal binding sites, each of which can bind a ferric ion with a $K_a$ of $\approx 10^{22}$ $M^{-1}$ together with a bicarbonate anion. The protein ligands for $Fe^{3+}$ at each site include two or three tyrosine residues, one or two histidine residues, and the concomitantly bound bicarbonate anion (1). The view that transferrin consists of two homologous domains—each associated with one metal binding site—is supported by the demonstration of internal homology in a partial sequence for human transferrin (2) and by the production of fragments of various transferrins by partial proteolysis that have approximately half the molecular weight of the native protein and single sites for $Fe^{3+}$ binding (3–6). Low resolution x-ray crystallographic studies also support a two-domain structure for rabbit serum transferrin (7).

The delivery of iron from transferrin to cells is mediated by the binding of transferrin–$Fe^{3+}$ complexes to specific cellular receptors (e.g., see refs. 8–10). Transferrin molecules therefore possess a specific receptor-recognition site in addition to the two metal binding sites.

We report here the complete amino acid sequence of human serum transferrin, which confirms the presence of extensive internal homology within the polypeptide chain and permits the identification of the locations of the substitutions in some pre-

viously reported genetic variants. Together with structural information from chicken ovotransferrin (11, 12), the locations of conserved residues of possible functional interest are identified.

## MATERIALS AND METHODS

Materials. Transferrin from pooled human sera was obtained from Behring (San Diego, CA) or from Sigma and was used without further purification. Trypsin (L-1-tosylamide-2-phenylethyl chloromethyl ketone-treated) and chymotrypsin were from Worthington; thermolysin from Calbiochem; Staphylococcus aureus protease from Miles; and pepsin from Sigma. Sequencer reagents were from Beckman or Pierce.

Methods. CNBr fragmentation of human transferrin was performed without prior cleavage of disulfide bonds. As reported previously (13), on separation by gel filtration with Sephadex G-75 in 5% formic acid, two groups of disulfide-bonded fragments (CNA and CNB) and three cystine-free fragments (CN7: residues 257–309, CN8: residues 383–389, and CN9: residues 310–313; see Fig. 2) were obtained. After reduction and alkylation with iodoacetamide or ethylenimine, fragments designated CN5 (residues 27–109) and CN6 (residues 1–26) were obtained from CNB by gel filtration, whereas CNA yielded five fragments: CN1 (residues 500–678), CN2 (residues 110–256), CN3 (residues 390–463), CN4 (residues 314–382), and CN10 (residues 465–499).

Digests of fragments with trypsin, chymotrypsin, thermolysin, or S. aureus protease were performed in 0.1 M ammonium bicarbonate or in unbuffered solution at pH 8.0 and 37°C for about 4 hr; protease to protein (wt/wt) ratios were between 3:100 and 5:100. Pepsin digests were performed in 5% formic acid. Cleavage at arginine residues was achieved by trypsin digestion of acetylated or citraconylated fragments (14). Peptides were purified by a variety of methods; in different cases, ion exchange chromatography with DEAE-cellulose (Whatman DE52) eluted with ammonium bicarbonate gradients, CM-cellulose (CM 52) eluted with gradients in pyridine acetate buffer at pH 5.0, or sulfonated polystyrene resins (Aminex A5) eluted with pyridine acetate gradients (15) were used. Impure fractions were repurified by gel filtration with Bio-Gel P4, Sephadex G-25 (superfine), or Sephadex G-50 (superfine) and by reverse-phase chromatography (16). Criteria for peptide purity were single NH2-terminal groups by the dansyl chloride procedure (17), stoichiometric ratios on amino acid analysis, and ultimately, sequence analysis. Amino acid analyses of acid hydrolysates were performed with a Durrum D-500 amino acid ana-

* Present address: Dept. of Biochemistry, Univ. of British Columbia, Vancouver, BC, Canada.
† Present address: Dept. of Endocrinology, Centro Ramon y Cajal, Madrid, Spain.
‡ Present address: Dept. of Biology, Massachusetts Institute of Technology, Cambridge, MA.
§ To whom reprint requests should be addressed.

lyzer. During earlier phases of this work, the sequences of peptides were determined by the dansyl chloride/Edman degradation procedure (18); for much of the sequence, automatic sequence analysis with a Beckman 890C sequencer was performed by employing the fast peptide dimethylallylamine program (Beckman program no. 021979) or the 0.1 M Quadrol program (program no. 121078) in combination with polybrene. Phenylthiohydantoin derivatives of amino acids were identified by gas chromatography and back hydrolysis with 6 M HCl containing 0.1% $SnCl_2$ (19) or—more recently—by HPLC (20).

## RESULTS AND DISCUSSION

The amino acid sequence of human transferrin was assembled by aligning the sequences of the 10 CNBr fragments. The structures of these fragments were separately determined by a combination of direct sequence analysis and from the structures of overlapping peptides obtained from a variety of proteolytic digests. The evidence for the sequences of the CNBr fragments, and some of the evidence for their order, is summarized in Fig. 1. Sequences for CN7, CN8, and CN9 have been published (18) and the evidence for their structures is omitted. The order of the fragments in the polypeptide chain was determined from a number of lines of evidence. (*i*) Peptides overlapping methionine residues were isolated from large fragments prepared by partial proteolysis of transferrin (unreduced) with pepsin. Three fragments isolated from such digests by gel filtration were cleaved by reduction and carboxamidomethylation and were reseparated. In some cases, further proteolysis with trypsin was performed, followed by reseparation. In this way peptides overlapping the following pairs of fragments were obtained: CN6→CN5, CN5→CN2, CN4→CN8, and CN3→CN10, together with a peptide containing the amino terminus of CN3 preceded by a methionine residue. (*ii*) From a digest of performic acid-oxidized transferrin with thermolysin, peptides overlapping CN7→CN9 and CN7→CN9→CN4 were obtained, as reported previously (2, 21). (*iii*) The locations of CN6 and CN1 at the $NH_2$- and COOH-terminal ends of the transferrin polypeptide chain are indicated by their structures. (*iv*) An Fe-binding fragment of $M_r$ 35,000, obtained by thermolysin digestion of diferric transferrin, contained CN6, CN5, CN2, and CN7 in a contiguous sequence, together with a fragment of CN4 (6). Therefore, CN7 must follow CN2 in the overall sequence.

Although overlaps were not obtained at every methionine residue, the combined evidence unambiguously indicates the order CN6→CN5→CN2→CN7→CN9→CN4→CN8→CN3→CN10→CN1. In numerous CNBr cleavages of transferrin, no fragments in addition to these 10 were observed, and amino acid analyses of transferrin are consistent with the presence of 9 methionine residues. Therefore, it can be concluded that the structure shown in Fig. 1 represents the entire amino acid sequence of transferrin.

**Disulfide Bond Arrangements.** The pairing of half-cystine residues has been examined by the isolation and analysis of cystine-containing peptides from a tryptic digest (pH 6.0) of CNB (CN6 and CN5). Peptides generated by reduction and carboxymethylation of these cystine peptides indicated that Cys-9 is linked to Cys-48 and Cys-19 to Cys-39. Cleavage of the thermolysin-produced Fe-binding "half molecule" with CNBr gave a disulfide-bonded fragment containing CN2 (residues 110–256) and part of CN4 (residues 326–341). Cystine-containing peptides isolated from a peptide digest of this fragment indicated the presence of disulfide bonds between half-cystines 117 and 194, cystines 137 and 331, and cystines 227 and 241. Half-cystines 158, 161, 171, 174, 177, and 179 were present in a single peptide that has proved difficult to subcleave. The arrangement

of these three bonds remains uncertain. Cys-339—which is also present in the thermolysin fragment—was found to be linked to a cystine in the sequence Ala-Cys, probably corresponding to residues 594–595. The arrangements of the remaining disulfide bonds in the carboxylterminal half of the molecule must still be investigated.

**General Features of the Molecule.** The polypeptide chain of human transferrin contains 678 amino acids and—together with the two glycan moieties, $M_r$ 2207 each (22)—has an overall molecular weight of 79,550. The most striking feature of the sequence is the presence of strong 2-fold internal homology previously suggested from partial sequence studies (2). As shown in Fig. 2, when residues 1–336 are aligned with residues 337–678 by the inclusion of gaps (23 in the former sequence and 19 in the latter) to optimize the similarity, 143 residues in corresponding positions (40%) in the two sequences are identical and a considerable proportion of the residues that are not identical are similar in chemical nature. The most reasonable hypothesis compatible with this finding is that the structural gene for the transferrin molecule arose during the course of evolution by the contiguous duplication of the structural gene for an ancestral protein of ≈340 amino acids. Together with previous studies of Fe-binding fragments of transferrin, the internal homology indicates that transferrin consists of two Fe-binding domains—corresponding to residues 1–336 and 337–678—and that the smaller ancestor of the modern transferrins was a single-domain molecule containing a single Fe-binding site. Studies of the disulfide bond arrangement show that Cys-331 is linked to a residue within the $NH_2$-terminal domain and Cys-339 to a residue deep within the COOH-terminal domain. Therefore, the two domains must be closely packed together with any linking region being limited to a maximum of 7 residues (332–338).

The functional significance of the presence of two domains with separate Fe-binding sites is uncertain as, although the two sites have some distinguishable physical properties (see ref. 1), present evidence indicates that in human transferrin there is no difference in the *in vivo* behavior of the sites with respect to iron uptake and delivery to cells (24). The evolutionary advantage of the doubled structure may instead lie in the reduction of losses on glomerular filtration.

The sequence reported here was determined with transferrin from pooled human sera and represents the structure of the predominant genetic form, transferrin C. No sign of sequence variability was observed in any regions of the structure. However, the sites of substitutions in the sequence in a number of genetic variants can be assigned on the basis of difference peptides isolated by previous workers. A chymotryptic peptide isolated by Wang and Sutton (25) from transferrin C appears to correspond to residues 275–282; the corresponding peptide from the D1 variant indicates a substitution of Gly for Asp at position 277.

Peptides isolated from a tryptic digest of CN7 of transferrin $D_{Chi}$ by Wang *et al.* (26) and Howard *et al.* (23) show that His-300 is replaced by Arg in this variant. On the basis of peptides isolated by Wang and co-workers (27), Gly-651 appears to be replaced by Glu in transferrin B2.

**Comparisons of the Domains.** When the nature of the amino acids conserved between the two domains is analyzed, a considerable proportion appears to be of potential structural significance (e.g., cystines, glycines, and hydrophobic residues)—possibly reflecting the preservation of a similar three-dimensional structure in the two domains. Some conserved residues can be expected to have a functional significance, as it can be reasonably expected that the stereochemical restrictions on the arrangements of liganding groups in the Fe-binding sites would result in a conservation of the residues involved in these

1 V P D K T V R W C A V S E H E A T K C Q S F R D H M K S V I P S D G P S V A C V K K A S Y L D C I R A I A A N E A D A V

61 T L D A G L V Y D A Y L A P N N L K P V V A E F Y G S K E D P E T F Y Y A V A V V K K D S G F Q M N Q L R G K K S C H T

121 G L G R S A G W N I P I G L L Y C D L P E P R K P L E K A V A N F F S G S C A P C A D G T D F P Q L C Q L C P G C G C S

181 T L D E Y F G Y S G A F K C L K D G A G D V A F V K H S T I F E N L A N K A D R D Q Y E L L C L D N T R K P V Q D Y K D

241 C H L A E V P S H T V V A R S M G G K E D L I W E L L N Q A Q E H F G K D K S K E F Q L F S S P H G K N L L F K D S A H

301 G F L K V P P R M N A K M Y L G Y E Y V T A I R N L R E G T C P E A P T N E C K P V K W C A L S H H E R L K C N E W S V

361 S D V G K I E C V S A E T T E D C I A K I M N G E A D A M S L D G G F V Y I A G K C G L V P V L A E N Y N K S D D C E Q

421 T P A D G Y F A V A V V K K S A S D L T W D N L K G K K S C H T A V G R T A G W N I P M G L L Y N K I N H C R F D E F F

481 S E G C A P G S K K D S S L C K L C M G S G L N L C E P N N K E G Y Y G Y T G A F R C L V E K G D V A F V K H Q T V T Q

541 N P G G K N P D W P A K D L N E K Y N E L C L D G T R K P V Q E Y A N C H L A R A P N H A V V T R K D K E A C V H K I L

601 R Q Q Q H L F G S N V T D C S G N F C L F R S E T K D L L F R D D T V C L A K L H D R N T Y E K Y L G Q E Y V K A V G N

661 L R K C S T S S L L E A C T F R R P

FIG. 1. The amino acid sequence of human serum transferrin. Thick bars indicate portions of fragments that were subjected to sequence analysis and thin bars, portions that were not. The highest bar below each line of sequence shows the region of each CNBr fragment that was determined by direct automatic sequence analysis. Peptides were generated as follows: T, trypsin (on AE, CAM, or CM fragments); C, chymotrypsin; TH, thermolysin; SP, *S. aureus* protease; PP, partial peptic hydrolysis of whole protein; PP/T, partial peptic hydrolysis with subsequent trypsin cleavage; Pf TH, thermolysin cleavage of performic-oxidized protein. Glycosylated asparagine residues are indicated by ★. The evidence for the structure of CN7 (residues 257–309) has been omitted as it was published previously (18).

binding sites. There are five histidine and nine tyrosine residues that are conserved in the two domains of human transferrin. Some of these can be eliminated from consideration as liganding residues as they are not conserved in corresponding sequences from chicken ovotransferrin. Thus, His-14 is changed to Pro (see ref. 28) and His-242 to Asp in a peptide that is clearly homologous with residues 237–249 of human transferrin [peptide 29, Elleman and Williams (11)]. In contrast, histidines 249, 119, and 451 are conserved [peptides 30 and 31, Elleman and Williams (11)]. Therefore, one or two of the homologous pairs of histidine residues 119/451, 207/535, or 249/584 must probably provide the imidazole ligands for binding the $Fe^{3+}$ ions. Of the conserved tyrosine residues (9 pairs), only the pairs at positions 314/649 and 319/654 can be eliminated, as residues 314 and 319 are components of the region missing from the amino terminal Fe-binding fragment obtained by thermolysin digestion (6).

There are some pronounced differences between the $NH_2$-terminal (N) and COOH-terminal (C) domains. Thus, the N-domain contains fewer disulfide bonds than the C-domain (8

```
1                    10                        20                        30
VAL PRO ASP [LYS] THR [VAL] ARG TRP CYS ALA VAL SER GLU [HIS GLU] ALA THR [LYS CYS] GLN SER PHE ARG ASP HIS MET LYS SER VAL [ILE]
ASN GLU CYS [LYS] PRO [VAL] LYS TRP CYS ALA LEU SER HIS [HIS GLU] ARG LEU [LYS CYS] ASN GLU TRP SER VAL SER ASP VAL GLY LYS [ILE]
337      340                   350                            360

                              40                  50                        60
PRO SER ASP GLY PRO SER VAL ALA [CYS VAL] LYS LYS ALA SER TYR LEU [ASP CYS ILE] ARG ALA [ILE] ALA ALA ASN [GLU ALA ASP ALA] VAL
GLU ————————————————————————————— [CYS VAL] SER ALA GLU THR THR GLU [ASP CYS ILE] ALA LYS [ILE] MET ASN GLY [GLU ALA ASP ALA] MET
                              370                      380

                              70                  80                        90
THR [LEU ASP] ALA [GLY] LEU [VAL TYR] ASP [ALA] TYR LEU ALA PRO ASN ASN [LEU] LYS [PRO VAL] VAL [ALA GLU] PHE [TYR] GLY SER LYS GLU [ASP]
SER [LEU ASP] GLY [GLY] PHE [VAL TYR] ILE [ALA] ——— GLY LYS CYS GLY [LEU] LEU [PRO VAL] LEU [ALA GLU] ASN [TYR] ASN LYS SER ASP [ASP]
                              390                 400                       410

                              [PRO] GLU THR PHE [TYR TYR] ALA VAL ALA VAL VAL LYS LYS —— ASP [SER] GLY PHE GLN MET ASN GLN [LEU] ARG [GLY LYS]
CYS GLU GLN THR [PRO] ALA ASP GLY [TYR] [PHE] ALA VAL ALA VAL VAL LYS LYS SER ALA [SER] ASP LEU THR TRP ASP ASN [LEU] LYS [GLY LYS]
                   420                          430                       440

                              120                      130                       140
LYS SER CYS HIS THR [GLY LEU] [GLY ARG] SER [ALA GLY TRP ASN ILE PRO] ILE [GLY LEU LEU TYR] CYS ASP LEU PRO GLU PRO [ARG] LYS PRO
LYS SER CYS HIS THR [GLY ALA VAL] [GLY ARG] THR [ALA GLY TRP ASN ILE PRO] MET [GLY LEU LEU TYR] ASN LYS ILE ASN HIS CYS [ARG] PHE ASP
                   450                          460                       470

                              150                      160                       170
LEU GLU LYS ALA VAL ALA ASN [PHE PHE SER] GLY SER [CYS ALA PRO] CYS ALA ASP GLY THR ASP PHE PRO GLN [LEU CYS] GLN [LEU CYS] PRO
————————————————————————————— GLU [PHE PHE SER] GLU GLY [CYS ALA PRO] ——— GLY SER LYS LYS ASP SER SER [LEU CYS] LYS [LEU CYS] MET
                   480                                   490

[GLY] CYS [GLY] ————————— [CYS] SER THR LEU ASP ———— GLU [TYR] PHE [GLY TYR] SER [GLY ALA PHE] LYS [CYS LEU] LYS ASP GLY ALA [GLY]
[GLY] SER [GLY] LEU ASN LEU [CYS] GLU PRO ASN ASN LYS GLU GLY [TYR] TYR [GLY TYR] THR [GLY ALA PHE] ARG [CYS LEU] —— VAL GLU LYS [GLY]
500                           510                       520

                              210                      220
ASP VAL ALA PHE VAL LYS HIS SER [THR] ILE PHE GLU [ASN] LEU ALA ASN [LYS] ———————————————— ALA ASP ARG ASP GLN [TYR]
ASP VAL ALA PHE VAL LYS HIS GLN [THR] VAL THR GLN [ASN] PRO GLY GLY [LYS] ASN PRO ASP TRP PRO ALA LYS ASP LEU ASN GLU LYS [TYR]
530                           540                       550

GLU [LEU] [LEU CYS LEU ASP] ASN THR ARG LYS PRO VAL [GLN] ASP [TYR] LYS ASP [CYS HIS LEU ALA] GLU VAL [PRO] SER HIS THR [VAL VAL] ALA
ASN [GLU] [LEU CYS LEU ASP] GLY THR ARG LYS PRO VAL [GLN] GLU [TYR] ALA ASN [CYS HIS LEU ALA] ARG ALA [PRO] ASN HIS ALA [VAL VAL] THR
      560                              570                        580

[ARG] SER MET GLY GLY [LYS GLU] ASP LEU ILE TRP GLU LEU [LEU] ASN [GLN] ALA [GLN] GLU HIS [PHE GLY] LYS ASP LYS SER LYS GLU PHE GLN
[ARG] LYS ASP ———— [LYS GLU] ALA CYS VAL HIS LYS ILE [LEU] ARG [GLN] GLN [GLN] HIS LEU [PHE GLY] SER ASN VAL THR ASP CYS SER GLY
  590                       600                       610

                              290                      300                       310
LEU PHE ———— SER SER PRO HIS GLY [LYS] ASN [LEU LEU PHE] LYS [ASP] SER ALA HIS GLY PHE LEU [LYS] VAL PRO PRO [ARG] MET ASN ALA
ASN PHE CYS LEU PHE ARG SER GLU THR [LYS] ASP [LEU LEU PHE] ARG [ASP] ASP THR VAL CYS LEU ALA [LYS] LEU HIS ASP [ARG] ASN THR TYR
      620                            630                          640

                              320                      330                       336
LYS MET [TYR LEU GLY] TYR [GLU] [GLU TYR VAL] THR [ALA] ILE ARG [ASN LEU ARG] GLU GLY THR CYS ———— PRO [GLU ALA] PRO [THR] ————
GLU LYS [TYR LEU GLY] GLN [GLU] [GLU TYR VAL] LYS [ALA] VAL GLY [ASN LEU ARG] LYS CYS SER THR SER SER LEU LEU [GLU ALA] CYS [THR] PHE ARG
      650                            660                          670

ARG PRO
678
```
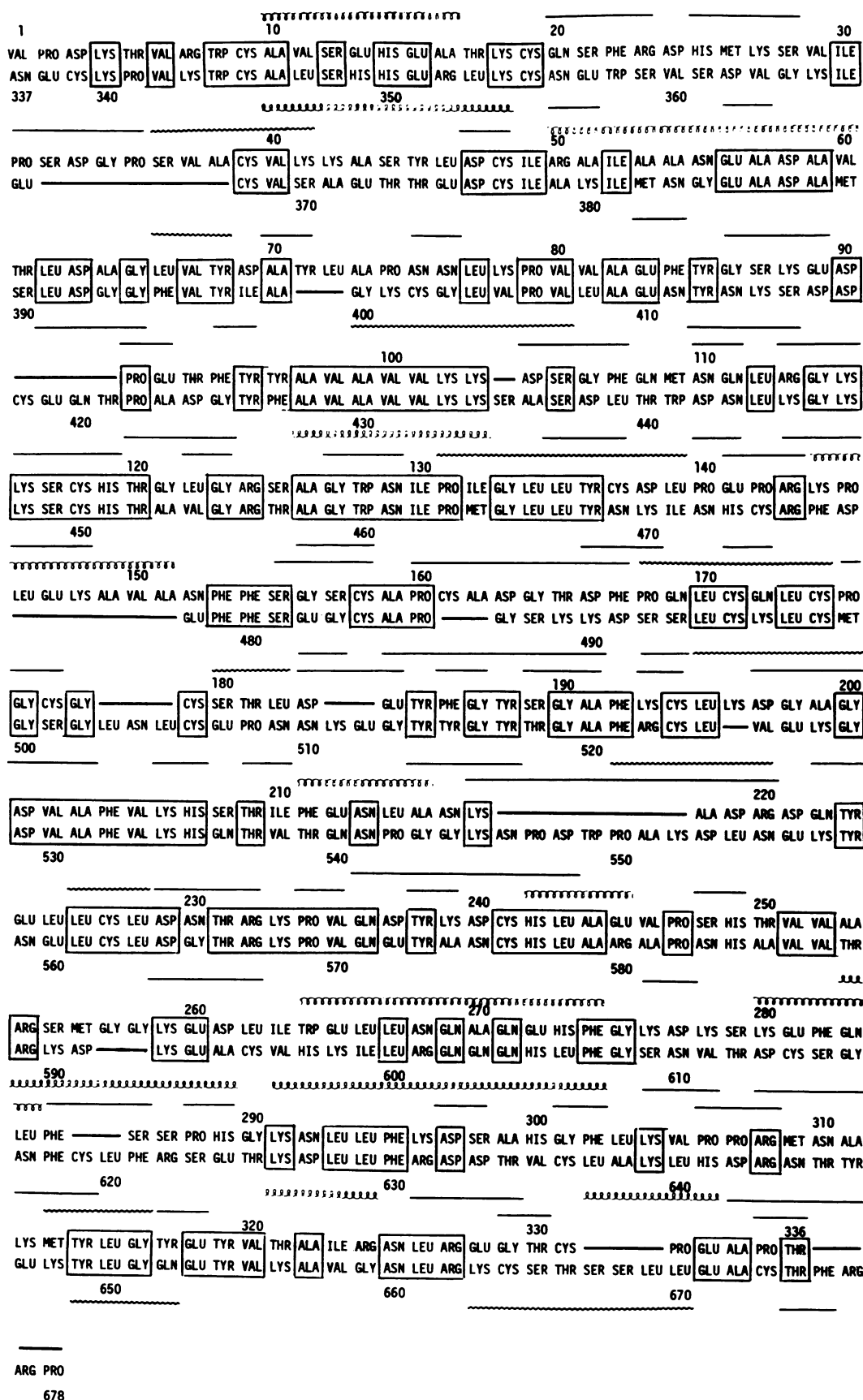
FIG. 2. Comparison of the primary structures and predicted secondary structures (23) of the two domains (residues 1–336 and 337–678). Gaps have been placed in the two sequences to maximize the homology. Boxes denote residues that are identical in the two domains. The symbols used for secondary structure features are ꝏꝏ, α-helix; ⌒⌒, extended conformation (β-structure); ——, bends.

compared to 11), whereas the C-domain contains both sites of glycosylation of the human transferrin molecule (asparaginine residues 411 and 610); the glycosylation sites are the only asparagines in human transferrin that are contained in N-glycosylation "signal sequences" of the type Asn-X-Ser/Thr. It is of some interest that the glycosylation site in chicken transferrin corresponds to residue 469 of human transferrin and is therefore also in the C-domain (12). The corresponding sequence in human transferrin, Asn-Lys-Ile, is not a potential site of glycosylation. This asymmetry of distribution of carbohydrate between the domains may be a general feature of transferrins (e.g., see ref. 5).

Many of the gaps in the alignment in Fig. 2—which presumably represent deletions that have occurred since the gene duplication that initially gave rise to the ancestral two-domain transferrin structure—are located in regions where the distribution of cystine residues differs in the two domains. These deletions may reflect structural adaptations to the addition or excision of disulfide bonds during the course of evolution.

**Secondary Structure.** Fig. 2 also shows regions of secondary structure predicted from the sequence by computer analysis (29). This indicates that 97 residues (14% of the polypeptide chain) have a high probability of being in an $\alpha$-helical conformation, whereas 70 residues (10%) are probably in an extended configuration. This is in reasonable agreement with the results of circular dichroism analyses of transferrin that suggest that 17% of the polypeptide chain is in an $\alpha$-helical conformation (30).

Many gaps in the sequence comparison between the domains are in regions that are predicted to have no ordered structure or to form a bend and may therefore be accommodated with little modification in the overall three-dimensional structure.

In contrast, the gap between residues 477 and 478 corresponds to 6 residues in the N-terminal domain that are predicted to be in an $\alpha$-helical conformation. In view of the considerable differences in structure between the domains between residues 137–184 and 469–513, this region may be the site of considerable conformational differences between the N- and C-domains. Although empirical predictions of secondary structure from amino acid sequences may be unreliable, the results are consistent with the view that the two domains have generally similar three-dimensional structures.

1. Aisen, P. & Listowsky, I. (1980) *Annu. Rev. Biochem.* **49**, 357–393.
2. MacGillivray, R. T. A., Mendez, E. & Brew, K. (1977) in *Proteins of Iron Metabolism,* eds. Brown, E. B., Aisen, P., Fielding, J. & Crichton, R. R. (Grune & Stratton, New York), pp. 133–142.
3. Williams, J. (1974) *Biochem. J.* **141**, 745–752.
4. Williams, J. (1975) *Biochem. J.* **149**, 237–244.
5. Brock, J. H. & Arzabe, F. R. (1976) *FEBS Lett.* **69**, 63–66.
6. Lineback-Zins, J. & Brew, K. (1980) *J. Biol. Chem.* **255**, 708–713.
7. Gorinsky, B., Horsburgh, C., Lindley, P. F., Moss, D. S., Parkar, M. & Watson, J. L. (1979) *Nature (London)* **281**, 157–158.
8. Speyer, B. E. & Fielding, J. (1974) *Biochim. Biophys. Acta.* **332**, 192–200.
9. van Bockxmeer, F., Hemmaplardh, D. & Morgan, E. H. (1975) in *Proteins of Iron Storage and Transport in Biochemistry and Medicine,* ed. Crichton, R. R. (Elsevier/North-Holland, Amsterdam), pp. 111–119.
10. Karin, M. & Mintz, B. (1981) *J. Biol. Chem.* **256**, 3245–3252.
11. Elleman, T. C. & Williams, J. (1970) *Biochem. J.* **116**, 515–535.
12. Kingston, I. B. & Williams, J. (1975) *Biochem. J.* **147**, 463–472.
13. Sutton, M. R. & Brew, K. (1974) *Biochem. J.* **139**, 163–168.
14. Dixon, H. B. F. & Perham, R. N. (1968) *Biochem. J.* **109**, 312–314.
15. Vanaman, T. C., Wakil, S. J. & Hill, R. L. (1968) *J. Biol. Chem.* **243**, 6420–6431.
16. Sinha, S. K. & Brew, K. (1981) *J. Biol. Chem.* **256**, 4193–4204.
17. Gray, W. R. & Hartley, B. S. (1963) *Biochem. J.* **89**, 379–380.
18. Sutton, M. R., MacGillivray, R. T. A. & Brew, K. (1975) *Eur. J. Biochem.* **51**, 43–48.
19. Mendez, E. & Lai, C. Y. (1975) *Anal. Biochem.* **68**, 47–53.
20. Hunkapiler, M. W. & Hood, L. E. (1978) *Biochemistry* **17**, 2124–2133.
21. MacGillivray, R. T. A. & Brew, K. (1975) *Science* **190**, 1306–1307.
22. Dorland, L., Hoverkamp, J., Schut, B. L., Vliegenart, J. F. G., Spik, G., Strecker, G., Foureit, B. & Montreuil, J. (1977) *FEBS Lett.* **77**, 15–20.
23. Howard, P. N., Wang, A.-C. & Sutton, H. E. (1968) *Biochem. Genet.* **2**, 265–269.
24. Huebers, H., Johnson, B., Huchers, E., Csiba, E. & Finch, C. (1981) *Proc. Natl. Acad. Sci. USA* **78**, 2572–2576.
25. Wang, A.-C. & Sutton, H. E. (1965) *Science* **149**, 435–437.
26. Wang, A.-C., Sutton, H. E. & Howard, P. N. (1967) *Biochem. Genet.* **1**, 55–59.
27. Wang, A.-C., Sutton, H. E. & Riggs, A. (1966) *Am. J. Hum. Genet.* **18**, 454–458.
28. Jolles, J., Mazurier, J., Boutique, M. H., Spik, G., Montreuil, J. & Jolles, P. (1976) *FEBS Lett.* **69**, 27–31.
29. Burgess, A. W., Ponnuswamy, P. K. & Scheraga, H. S. (1974) *Isr. J. Chem.* **12**, 239–286.
30. Nagy, B. & Lehrer, S. S. (1972) *Arch. Biochem. Biophys.* **148**, 27–36.