# PREDICTION REPORT

# A Consensus Prediction of the Secondary Structure for the 6-Phospho-β-D-Galactosidase Superfamily

Dietlind L. Gerloff and Steven A. Benner
*Department of Chemistry, Swiss Federal Institute of Technology, CH-8092 Zürich, Switzerland*

**ABSTRACT** Two separate unrefined models for the secondary structure of two subfamilies of the 6-phospho-β-D-galactosidase superfamily were independently constructed by examining patterns of variation and conservation within homologous protein sequences, assigning surface, interior, parsing, and active site residues to positions in the alignment, and identifying periodicities in these. A consensus model for the secondary structure of the entire superfamily was then built. The prediction tests the limits of an unrefined prediction made using this approach in a large protein with substantial functional and sequence divergence within the family. The protein belongs to the (α–β class), with the core β strands aligned parallel. The supersecondary structural elements that are readily identified in this model is a parallel β sheet built by strands C, D, and E, with helices 2 and 3 connecting strands (C + D) and (D + E), respectively, and an analogous β–α unit (strand G and helix 7) toward the end of the sequence. The resemblance of the supersecondary model to the tertiary structure formed by 8-fold α–β barrel proteins is almost certainly not coincidental. © 1995 Wiley-Liss, Inc.

Key words: prediction contests, α–β barrel, protein sequence alignment

A central problem in protein chemistry challenges the chemist to deduce the conformation (secondary and tertiary structure) of a protein from sequence information (primary structure). Both at the ETH in Zurich[1] and elsewhere,[2–6] progress toward solution of this problem has come through an analysis of patterns of conservation and variation in the sequences of homologous proteins.[7] Such an analysis is especially powerful when it is aided by detailed models of divergent evolution.[8] Predictions made using this approach are "consensus" models for conformation of a protein family, and assume that proteins related by common ancestry have similar conformations.[9]

The value of these methods can be explored by using them to make bona fide predictions, those published before an experimental structure becomes available. To date, over a dozen bona fide predictions have been made using these methods [reviewed in refs. 7 and 10]. For about half of these, a subsequently determined crystal structure has allowed these predictions to be evaluated. In many cases, the predictions have proven to be remarkably accurate.[10] It is now clear that predictions are possible that miss no core secondary structural elements, misassign no α helices as β strands (or vice versa), and do not overpredict any significant secondary structural element.[11] Predictions meeting these criterion are satisfactory as starting points to assemble a tertiary structural model of a protein family. Predicted secondary structures for pleckstrin homology domain,[12,13] hemorrhagic metalloproteinases,[14] and Src homology 2 domains[2,3] come close to meeting this standard.

Ongoing bona fide prediction efforts are necessary to define the scope of prediction methods. Over time, a large set of examples will emerge that will become statistically representative of proteins as a whole. As this set has accumulated to date, it has become clear that misassignments made by evolutionary analyses come in five principal types:

1. where multiple alignment is incorrect;
2. where the secondary structure of homologous proteins has diverged;
3. when attempting to distinguish between surface β strands and surface loops;
4. when attempting to distinguish between long internal β strands and internal helices;
5. when attempting to assign secondary structure to active site regions.

The first two (and often the first three) problems are interrelated. When the secondary structure has diverged, this often creates bad multiple alignments. Further, the distinction between a surface strand and a surface loop is often difficult even when the experimental data are in hand, and these elements

often undergo substantial divergence in conformation during divergent evolution.

A challenge was issued on October 10, 1994, to predict the secondary and elements of the supersecondary structure of the 6-phospho-$\beta$-D-galactosidase superfamily. The prediction was due before November 1. This protein family appeared to be an excellent target for placing the method to an extreme test. The protein is large; the target sequence has 468 amino acids. The family appears to adopt quaternary structure, at least in some cases. Both the thioglucosidase from *Brassica napus* (rape) and the thioglucosidase from *Sinapis alba* (white mustard) are reported to be homodimers, while the $\beta$-galactosidase from *Sulfolobus solfataricus* (not shown in the alignment) is a homotetramer. This implies that quaternary contacts might bury some residues that are on the surface of subunits, complicating the secondary structure prediction. Further, biological function has diverged substantially within the protein family, as measured by a wide divergence in substrate specificity in the member proteins. Finally, with only a few days to make a prediction, the example tests the ability of a prediction method to produce an accurate model without the benefit of extensive refinement.

A multiple alignment for the protein family was built from sequences extracted from SwissProt 29[15] using the DARWIN system.[16,17] Surface and interior residues were assigned by automated procedures similar to those described elsewhere,[18] the multiple alignment was parsed into units forming independent secondary structures, and elements of secondary structure were predicted within the parsed segments from patterns of conservation and variation, as described elsewhere.[10,12,14,19] Many of the automated routines used in this prediction are available to the public on a server accessible via electronic mail at the address cbrg@inf.ethz.ch, or using the World Wide Web (WWW) with URL http://cbrg.inf.ethz.ch/.

The secondary structure prediction is presented residue-by-residue in Figure 1, and summarized in Table 1. A summary of the secondary structure prediction follows:

Strand A (a009–a011; b049–051) is a short internal segment confirmed in both subfamilies.

Strand B (a014–020; b053–060) is separated from strand A by a GG dipeptide in a well anchored region of the alignment. It is largely internal. This region is interesting from a methodological point of view, as a strong assignment would not have been possible if only one of the two subfamilies were available. In both subfamilies considered alone, an internal helix would be possible. Together, however, a GG dipeptide parse at (a012–103) and a GG dipeptide parse at positions (b062–063), together with strong alignment anchoring excludes an internal helix in this region.

Parse region (a048–059; b071–081) is problematic. Subfamily b could contain a $\beta$ strand in this region (b073–078). However, it is matched with a parsing string (PGDSG; a050–054) in sequence e of subfamily a, and a strand was not assigned.

Helix 1 (a072–084; b095–107) is reliably assigned in both subfamilies, is well anchored, and displays good amphiphilicity.

Strand C (a089–093; b111–115) is problematic in subfamily a, in part because of the small number of sequences available in this subfamily. In subfamily b, the surface and interior assignments display alternating periodicity, which confirms the strand assignment.

Active site a (a095–102; b117–125) contains conserved Arg, Ser, Trp, and Arg. It is strongly assigned.

Helix 2 (a116–130; b138–153) is reliably assigned in both subfamilies, is well anchored, and displays good amphiphilicity.

Strand D (a136–140; b159–163) is well parsed in subfamily 2, is confirmed in both subfamilies, and is largely internal.

Active site b (a141; b164–166) contains conserved Thr (part of the preceding strand) and His.

Helix 3 (a158–177; b181–198) is reliably assigned in both subfamilies, is well anchored, and displays good amphiphilicity.

Strand E (a182–185; b205–208) is assigned to a region that is near the active site, where conservation associated with active site function often obscures patterns of variation and conservation that might be used to assign secondary structure.

Active site c (a184–187; b205–209) is relatively weak, and is based ultimately on a single conserved Asn. A conserved Thr two residues before supports this conclusion. Interesting, a Trp two residues earlier is almost completely conserved in the superfamily as well, as is a Glu immediately following the conserved Asn.

A region (b212–215) following this active site segment might be assigned as a $\beta$ strand in subfamily b. It is not paired with a reliable assignment in subfamily a, which contains repeated parsing elements that almost certainly exclude a standard secondary structural element in this area. A similar $\beta$ strand might be assigned in subfamily b (b219–221); this again has no corresponding element in subfamily a, and might form a $\beta$ hairpin with the preceding strand. The alignment is poorly anchored in this region, and considerable sequence divergence between the two subfamilies is evident.

Helix 4 (a212–226; b248–268) is cleanly amphiphilic up to position a227, when an interior assignment appears on the surface arc of the amphiphilic helix. The following segment also forms a short (8 residues) amphiphilic helical pattern. In subfamily b, the helix is largely internal. Nevertheless, to the extent that amphiphilicity is detected, it

extends past the position where the amphiphilic pattern is broken. This indicates that the contacts made in subfamily a are different from those made in subfamily b. Interestingly, this helix contains a conserved His (a218; b255) and a nearly conserved His (a224; b261).

Strand x (a242–245) is cleanly parsed in subfamily a, and is canonically assigned as a short β strand. The segment is disrupted by parsing elements in subfamily b, which appears to be well anchored. It is possible to identify a plausible β segment in this subfamily. Our experience, however, has been that the experimental assignments made for such regions depend strongly on the experimental secondary structure assignment tool.

Helix x (a259–273) is not cleanly amphiphilic (position a269), but is assigned nevertheless when considering subfamily a alone. A gap is placed in its middle in subfamily b (positions b311–312). If the multiple alignment of subfamily b2 is rearranged, a helix can be detected from positions (b303–317; total length 13 positions). If the multiple alignment of subfamily b1 is adjusted, and the sequence with the deletion discarded, a weak helix can also be found. The ambiguous alignment makes all of these assignments insecure, however, and there is significant possibility that the conformations of different members of the superfamily are quite different.

Strand y (a275–280) is assigned in subfamily a only. It corresponds to a parsed region in subfamily b. Two interior residues (b323–324) might form a corresponding structure, however, in subfamily b.

The amphiphilicity of helix 5 (a286–293; b332–342) is difficult to detect when examining the alignment overall. Examining subalignments, especially of subfamily b1 and subfamily b2, makes the amphiphilicity clearer.

The region (a314) might be assigned as a short helix (7–10 residues) if the left side of subfamily a is examined alone. There is no confirmation of this helix elsewhere, however, as this region of the alignment has undergone massive sequence divergence.

Strand F (a323–327; b381–388) is badly parsed in subfamily a. The segment is conceivably a continuation of a putative helix that may follow. In subfamily b, the strand is more reliably assigned. An excellent set of anchors aligns the subalignments, and we have chosen on these grounds to make the assignment definitive in the consensus secondary structure model.

Helix y (a329–339) is short, and contains a problematic residue at position (a336). There is no confirmation for a helix assignment in subalignment b. The ambiguous alignment makes this assignment further insecure, and there is a significant possibility that the conformations of different members of the superfamily are quite different.

Strand z (a375–382) is assigned in a region of the multiple alignment that has undergone massive sequence divergence, and where DARWIN had extreme difficulties achieving a plausible matching. It has plausible amphiphilicity in subfamily a. Therefore, the multiple alignment in subfamily b was collapsed in an effort to obtain regions that might also form β strands. For subfamily b1, segment (b446–452) displayed an alternating pattern. For subfamily b2, this was not possible, although it cannot be excluded that further rearrangement of the multiple alignment upon refinement could find an analogous region. As time was inadequate to do a complete search of different possible multiple alignments, no strand was assigned in this region in the consensus model.

Helix 6 (a385–398; b456–469) is well parsed, well anchored, amphiphilic, and confirmed in both subfamilies. It might, however, be missing one turn in some proteins in subfamily b.

Strand G (a404–407; b476–479) is well parsed, internal, and confirmed in both subfamilies.

Active site d (a408–410; b480–482), containing conserved Glu, Asn, and Gly, is not strongly assigned by analysis of the sequences themselves. It is, however, supported by biochemical work.[20]

Helix 7 (a431–448; b497–517) is well parsed, well anchored, amphiphilic, and confirmed in both subfamilies.

Residues (a451–a482; b522–554) form a remarkable segment. In subfamily b, the segment is not parsed for 35 residues, has a large number of interior residues, and apparently contains more than one secondary structural element. The first task is to parse this section. To this end, four additional columns were added to the multiple alignment by recognizing that lactase phlorizin hydrolase has multiple internal repeats. Interestingly, in two of these repeats, a parsing string PG appears. However, the repeats that contain this parsing string are cleaved proteolytically during the posttranslational modification.[1] These repeats are also missing Glu (b480), presumed to be part of an active site. Thus, there is no guarantee that these repeats have divergently evolved under functional constraints. This example makes an important point regarding the analysis of homologous sequences in the prediction of a protein structure.

In this region, an internal helix must be considered. Assignment of internal helices (as opposed to internal strands) relies on accurate parsing. The two subalignments were first carefully anchored. A reliable parse at (a471) was matched with a weak parse at (b541). A dipeptide GP parse in subfamily a (a460–462) was used to divide the first part of this segment. The conserved Asp was assumed to also indicate a break in secondary structure (as opposed to being an indicator of an active site position). This led to the assignment of four secondary structural elements in this region as follows:

Subfamily a

| Pos | cba def SIAPred | | Parse |
|---|---|---|---|
| 001 | M M M | M | |
| 002 | S T S | s | |
| 003 | KKK | N- | s |
| 004 | QTT | P- | I |
| 005 | LLL FFF | I | |
| 006 | PPP PPP | S | |
| 007 | QEK EAE | A | |
| 008 | DDD SHT | I | strand internal |
| 009 | FFF FFF | I | strand |
| 010 | VII LLL | I | strand |
| 011 | MFF WWW | . | strand |
| 012 | GGG GGG | - | parse |
| 013 | GGG GGG | s | parse |
| 014 | AAA AAA | I | strand internal |
| 015 | TTT LIT | i | strand |
| 016 | AAA AAA | I | strand |
| 017 | AAA AAA | i | strand |
| 018 | YYY NNN | I | strand |
| 019 | QQQ QQQ | i | strand |
| 020 | VAA SVV | A | strand |
| 021 | EEE EEE | i | active site? |
| 022 | GGG GGG | I | parse |
| 023 | AAA AAA | i | parse |
| 024 | TTT FYW | I | parse |
| 025 | KNH RLQ | S | parse |
| 026 | ETT ETE | s | parse |
| 027 | DDD GDD | s | parse |
| 028 | GGG DGE | A | parse |
| 029 | KKK KKK | i | |
| 030 | GGG GGG | i | |
| 031 | LLI | i | |
| 032 | TSS | s | |
| 033 | TTT | a | |
| 034 | VSS | | |
| 035 | DDD | a | |
| 036 | MLL | i | |
| 037 | IQQ | s | |
| 038 | PPP | s | |
| 039 | HQH | | |
| 040 | GGG | - | |
| 041 | _IV | i | |
| 042 | _FM | - | |
| 043 | _GG | s | |
| 044 | _EK | s | |
| 045 | _IM | i | |
| 046 | EVE | - | |
| 047 | HTP | a | |
| 048 | RRE RRR | s | parse |
| 049 | VVV MQI | i | |
| 050 | LAA AEL | i | parse |
| 051 | WWW VGG | s | parse |
| 052 | DDD KQK | s | parse |
| 053 | DYK LSE | a | parse |
| 054 | FYY GGN | s | parse |
| 055 | LLL LLI | i | |
| 056 | DEE EKK | S | |
| 057 | KED K_ | S | |
| 058 | QNN R_ | s | |
| 059 | GYY P_ | i | |
| 060 | ___Q | . | |
| 061 | _L R_ | . | |
| 062 | ___D | . | |
| 063 | ___D | . | |
| 064 | ___E | . | |
| 065 | RWW P_ | i | |
| 066 | FYY Y_ | s | |
| 067 | KTT P_ | i | |
| 068 | PAA S_ | . | |
| 069 | DEE HDD | s | parse |
| 070 | EPP EVV | I | parse |
| 071 | PPP EVV | i | |
| 072 | AAA AAA | I | helix amphiphilic |

Subfamily b

| Pos | e ba dc f lj ig kh SIAPred | | Parse |
|---|---|---|---|
| 040 | L LLL | i | |
| 041 | R SS | s | |
| 042 | S SS | s | |
| 043 | S KK | i | |
| 044 | S NN | - | |
| 045 | F FF FF | I | |
| =046 | E PE EE PE | S | |
| -047 | A KK SQ | S | parse |
| -048 | T GD DD | S | parse |
| =049 | F FF FF | I | parse |
| -050 | M LI KM | L | strand internal |
| -051 | W WW WW | F | strand |
| -052 | G GG GG | S | strand |
| 053 | T AS VT | V | parse |
| -054 | A GA AS | S | strand internal |
| 055 | T TT TT | T | strand |
| -056 | S AA AA | A | strand |
| -057 | Y SA AA | S | strand |
| 058 | Y EY YY | F | strand |
| -059 | I QQ QQ | Q | strand |
| 060 | I II II | Y | strand |
| -061 | E EE EE | E | active site? |
| =062 | G AA AS | A | parse |
| -063 | G GA AS | S | parse |
| -064 | I WY YT | N | parse |
| 065 | D NN NQ | K | parse |
| 066 | E EE EE | A | parse |
| -067 | D DD DD | N | parse |
| -068 | G DD DD | N | parse |
| -069 | R KR KK | R | parse |
| -070 | T GG GG | G | parse |
| 071 | P EE ML P | P VV LL | I | |
| 072 | E SS NN SG | S NN SG | s | |
| 073 | I II VV II | I VV II | I | |
| =074 | D DD DD DD | W WW WW | A | |
| =075 | G DD DD DD | N WW WW | s | |
| 076 | T RR FF FF | F GG FF | s | |
| -077 | F TS SS SS | A TS SS | I | |
| 078 | C TS SS AA | C SS AA | i | |
| 079 | D HH HH HH | N HH HH | I | |
| 080 | I QT TT TM | M TT TM | s | |
| 081 | I KE EE EE | E EE EE | s | |
| 082 | G RG GG GR | G GG GR | i | |
| 083 | G NK KK KR | N KK KR | s | |
| 084 | K HH KK RK | V KK RK | s | |
| 085 | V LA KF P | A LA KF | a | |
| 086 | Y II II VV | Y II VV | i | |
| 087 | D YD NN G | D NN G | s | |
| 088 | D DD DD DD | K DD DD | s | |
| 089 | G GG GG GR | H GG GR | i | |
| 090 | C HH DD HH | N DD HH | . | |
| 091 | G GG GG GG | D GG GG | . | |
| 092 | G GG GG GG | A GG GG | . | |
| 093 | V VV VV VV | I VV VV | i | |
| 094 | V TT VV | I TT VV | I | |
| -095 | A AA AA AA | A AT ST AA | S | helix amphiphilic |

| Pos | | | Parse |
|---|---|---|---|
| 073 | ASS TII | i | helix |
| 074 | DDD DDD | A | helix |
| 075 | FFF FFF | I | helix |
| 076 | HHH HHH | I | helix |
| 077 | HHH HHH | S | helix |
| 078 | RRK RRR | S | helix |
| 079 | YYY YYY | s | helix, parse? |
| 080 | DPP KPP | s | helix, parse? |
| 081 | EVV EQE | A | active site? |
| 082 | DDD DDD | I | helix |
| 083 | LLL III | I | helix |
| 084 | AEE AAA | I | helix |
| 085 | LLL LLL | i | break to the inside |
| 086 | ASA MFP | i | |
| 087 | EEE AAA | S | strand? |
| 088 | KKE EEE | I | strand? |
| 089 | YFY MMM | i | strand? |
| 090 | GGG GGG | i | strand? |
| 091 | HVV FFF | s | strand? |
| 092 | QNN KTT | i | strand? |
| 093 | VGG VCC | I | active site |
| 094 | III FLL | i | |
| 095 | RRR RRR | I | active site |
| 096 | VII TTI | i | |
| 097 | SSS SSS | i | |
| 098 | III III | I | |
| 099 | AAA AAA | I | active site |
| 100 | WWW WWW | I | |
| 101 | SSS SFA | i | active site |
| 102 | RRR RRR | A | active site |
| 103 | III LII | I | |
| 104 | FFF FFF | I | |
| 105 | EPP PPP | s | parse |
| 106 | DNT QQQ | . | parse |
| 107 | GGG GGG | i | parse |
| 108 | AXY | . | |
| 109 | GGG DDD | A | |
| 110 | EEE EEE | I | |
| 111 | VVV IAV | I | |
| 112 | ENN TEE | S | parse |
| 113 | PPE PPE | a | parse |
| 114 | __ NNN | s | parse |
| 115 | __ OEE | . | |
| 116 | RKK QAA | s | helix amphiphilic |
| 117 | GGG GGG | I | helix |
| 118 | VVV ILL | . | helix |
| 119 | AEE AAA | I | helix |
| 120 | FYF FFF | i | helix |
| 121 | YYY YYY | I | helix |
| 122 | HHH RDD | s | helix |
| 123 | KKK SRR | S | helix |
| 124 | LLL VLL | I | helix |
| 125 | FFF FFF | I | helix |
| 126 | AAA EDD | i | helix |
| 127 | DEE EEE | S | helix |
| 128 | CCC CLM | I | helix |
| 129 | AHH KAA | s | helix |
| 130 | AKK FKQ | . | break to surface |
| 131 | HRR YYA | I | |
| 132 | HHH GGG | I | |
| 133 | IVV III | . | |
| 134 | EEE EQK | S | parse? |
| 135 | PPP PPP | I | strand |
| 136 | FFF LLL | I | strand |
| 137 | VVV VVV | A | strand |
| 138 | TTT TTT | I | strand |
| 139 | LLL LLL | . | strand |
| 140 | HHH CSS | . | active site |
| 141 | HHH HHH | A | |
| 142 | FFF FYY | I | active site |
| 143 | DDD DEE | . | parse |
| 144 | TTT VMM | . | |
| 145 | PPP PPP | s | surface strand? |
| 146 | EEE MYY | I | |
| 147 | RVA HGG | . | parse |
| 148 | LLL LLL | s | |
| 149 | HHH VVV | i | |
| 150 | EKS TEK | S | parse |

| Pos | | | |
|---|---|---|---|
| 096 | C CC C | C | I | helix |
| -097 | D DD DD | D | s | helix |
| 098 | H HH SS | H | I | helix |
| -099 | F YY YY | Y | s | helix |
| -100 | H HH HH | I | helix |
| -101 | R RR RR | R | I | helix |
| -102 | F FY EE | K | S | helix |
| 103 | E EE EE | E | A | helix active site |
| 104 | D DD DD | D | I | helix |
| =105 | D DD DD | D | S | helix |
| -106 | V VI VL | L | I | helix |
| 107 | Q SK QR | D | I | break to inside |
| -108 | I IV TA | L | i | |
| 109 | M MM LL | L | s | strand |
| 110 | K KG QQ | s | strand |
| 111 | Q LI LL | M | s | strand |
| 112 | I DE EE | N | s | strand |
| -113 | G GG GG | G | . | strand |
| 114 | E QQ EE | . | active site |
| 115 | L KK KR | A | . | active site |
| 116 | H AS VT | I | i | active site |
| 117 | Y RR RR | A | i | active site |
| -118 | F FF FF | F | i | |
| 119 | F FF FF | I | active site |
| -120 | V II FF | L | i | |
| -121 | V II SS | I | I | active site |
| -122 | A AS AS | A | i | |
| -123 | W WW WW | W | I | active site |
| -124 | T PP PP | P | S | active site |
| -125 | R RR RR | R | i | |
| 126 | I II VT | L | i | parse |
| -127 | M FF LF | F | . | parse |
| -128 | A DE QE | N | s | parse |
| -129 | A DE QE | Q | s | parse |
| -130 | - GG GG | G | -i | parse |
| 131 | - LL VV | . | parse |
| 132 | A FT TF | F | s | parse |
| 133 | A FT TF | G | i | |
| 134 | G GG RR | N | a | |
| -135 | I TK EE | Y | s | break |
| -136 | I VL NN | I | S | helix amphiphilic |
| -137 | N NN NN | N | S | helix |
| -138 | E QQ RO | E | s | helix |
| -139 | E EE AQ | A | I | helix |
| =140 | G GG GG | G | I | helix |
| -141 | L ED DD | F | . | helix |
| -142 | Y FF YY | Y | I | helix |
| -143 | Y YY YY | D | I | helix |
| =144 | E DK HH | D | s | helix |
| =145 | H RR RR | R | i | helix |
| -146 | R NN NK | K | I | helix |
| -147 | L IT VV | V | S | helix |
| 148 | I LV LV | I | I | helix |
| 149 | D NN DD | DD | A | helix |
| 150 | E KL LL | C | S | helix |
| 151 | E VL LL | i | helix |
| -152 | E AQ AA | A | S | helix |
| 153 | E AN NN | K | A | helix |
| 154 | A NN KK | R | I | break |
| 155 | G SG SG | K | s | |
| 156 | L II II | K | S | |
| =157 | I EM PP | . | parse? |
| 158 | P PP PP | QK | S | parse? |
| 159 | M VA PF | Y | i | strand |
| -160 | L VT CC | A | I | strand |
| 161 | T TT TT | T | I | strand |
| -162 | I LL TM | A | I | strand |
| -163 | L YY YY | H | A | strand |
| =164 | W WW WF | H | I | active site |
| 165 | W WW WW | H | i | |
| -166 | D DD DL | L | s | parse |
| 167 | L LL LL | P | I | |
| 168 | P PP PP | P | . | surface strand? |
| 169 | Q WK AA | L | s | |
| -170 | W KA II | H | i | |
| -171 | I FF YY | H | . | parse |
| 172 | E LL QQ | Q | s | |
| 173 | D DD DD | D | s | |

Protein sequence alignment with secondary-structure annotations.

| # | Sequence | | Indicator | Annotation |
|---|---|---|---|---|
| 151 | ADN | EKN | S | parse |
| 152 | | YHY | s | parse |
| 153 | GGG | GGG | . | parse |
| 154 | DDD | SGG | s | parse |
| 155 | WFF | WWW | s | I |
| 156 | LLL | RGA | s | helix |
| 157 | SNN | NNN | s | helix |
| 158 | QRR | RRR | S | helix |
| 159 | EKE | KLA | I | helix |
| 160 | MTN | LTV | S | helix |
| 161 | LLI | VII | I | helix |
| 162 | DDE | EDD | s | helix |
| 163 | DYH | FCH | s | helix |
| 164 | FFF | FFF | I | helix |
| 165 | VVI | SEE | S | helix |
| 166 | ADD | RRH | I | helix |
| 167 | YYY | YYY | I | helix |
| 168 | AAA | AAA | I | helix |
| 169 | KEA | RRR | i | helix |
| 170 | FYF | TTT | I | helix |
| 171 | CCC | CVV | I | helix |
| 172 | FYF | FFF | I | helix |
| 173 | EKE | EAT | S | helix |
| 174 | EEE | ARR | s | helix |
| 175 | FFF | FYY | I | helix |
| 176 | SPP | DRQ | s | helix |
| 177 | EEE | GHH | s | helix |
| 178 | | LKK | I | |
| 179 | VVV | VVV | I | |
| 180 | KKN | KKA | s | |
| 181 | YYY | YRL | I | strand |
| 182 | WWW | WWW | I | strand |
| 183 | IFT | LLL | A | strand |
| 184 | TTT | TTT | I | strand |
| 185 | IFF | FFF | A | active site |
| 186 | NNN | NNN | A | |
| 187 | EEF | EEE | I | |
| 188 | PII | III | s | parse |
| 189 | TGG | NNN | s | parse |
| 190 | SPE | IMM | i | parse |
| 191 | MII | MSS | I | |
| 192 | AGG | LLL | I | |
| 193 | VDD | HHH | I | |
| 194 | QGG | SAA | s | |
| 195 | QQQ | PPP | I | parse |
| 196 | AGG | EEE | I | parse |
| 197 | TLL | STT | s | parse |
| 198 | GGG | GGG | s | parse |
| 199 | YYY | AVV | s | parse |
| 200 | GGG | GGG | B | |
| 201 | TKK | LLL | | |
| 202 | FFF | VPA | | |
| 203 | PPP | PPE | | |
| 204 | PPP | SAA | | |
| 205 | AGG | ESS | | |
| 206 | EII | GDD | | |
| 207 | SKK | EKE | | |
| 208 | GYY | NAA | | |
| 209 | RDD | Q | | |
| 210 | FPL | D | | |
| 211 | DBA | Q | | |
| 212 | KKK | VAE | | |
| 213 | TVV | KIV | | |

| # | Sequence | | Indicator | Annotation |
|---|---|---|---|---|
| 214 | FFF | YYY | | |
| 215 | QQQ | QQQ | | |
| 216 | ASS | AAA | | |
| 217 | EHH | AII | | |
| 218 | HHH | HHH | | |
| 219 | NNN | HHH | | |
| 220 | QMM | QQQ | | |
| 221 | MMM | LLL | I | helix |
| 222 | VVV | VVV | A | helix |
| 223 | AAS | AAA | s | helix |
| 224 | HHH | SSS | I | helix |
| 225 | AAA | AAA | I | helix |
| 226 | RRR | LRR | S | helix |
| 227 | TAA | AAA | I | helix |
| 228 | VVV | TVV | I | helix |
| 229 | NKK | KKK | s | helix |
| 230 | LLL | IAA | I | helix |
| 231 | YFY | ACC | I | helix |
| 232 | KKK | HHH | I | helix |
| 233 | SDD | EDS | I | helix (break to inside) |
| 234 | MGK | VML | i | helix |
| 235 | QGG | NIL | s | parse |
| 236 | LYY | PPP | i | parse |
| 237 | GKK | QDE | S | parse |
| 238 | GGG | NAA | s | parse |
| 239 | QEE | QQK | s | parse |
| 240 | III | VII | I | |
| 241 | GGG | GGG | S | |
| 242 | IVV | ONN | I | strand |
| 243 | VVV | MMM | I | strand |
| 244 | HHH | LLL | I | strand |
| 245 | AAA | ALL | I | strand |
| 246 | LLL | GGG | A | active site |
| 247 | QPP | GAG | s | parse |
| 248 | TTT | NML | s | parse |
| 249 | VKK | PLV | i | parse |
| 250 | YYY | YYY | I | strand, internal |
| 251 | PPP | PPP | I | parse? |
| 252 | YPY | YLL | I | strand |
| 253 | SDD | STT | s | strand |
| 254 | DPE | CSC | S | strand |
| 255 | SSE | KKQ | S | parse |
| 256 | ANN | PPP | . | parse |
| 257 | VPE | | I | |
| 258 | EA | | s | s |
| 259 | DDD | EEQ | s | helix |
| 260 | HVV | DDD | s | helix |
| 261 | HRR | VVM | . | si |
| 262 | AAA | WML | I | helix |
| 263 | EEE | ASA | S | helix |
| 264 | LLL | LLM | I | helix |
| 265 | QEE | EHE | S | helix |
| 266 | DDD | KQE | s | helix |
| 267 | AII | DNN | I | helix |
| 268 | LII | RRR | S | helix |
| 269 | EHH | EER | s | helix |
| 270 | NNN | NWW | I | helix |
| 271 | RKK | LLM | I | helix |
| 272 | LFF | FFF | I | break to inside |
| 273 | YII | FFF | i | strand |
| 274 | LLL | IGG | i | strand |
| 275 | DDD | DDD | i | strand |
| 276 | GAA | VVV | i | strand |
| 277 | TTT | QQQ | i | strand |
| 278 | LYY | AVA | . | strand |
| 279 | ALL | RRR | | |
| 280 | GGG | GGG | s | parse |
| 281 | EKH | TRQ | | parse |

(far right column)

| # | Sequence | | Indicator | Annotation |
|---|---|---|---|---|
| 251 | F ML | YY | i | helix |
| 252 | T DE | II | si | helix |
| 253 | A VV | VV | I | helix |
| 254 | A VS GG | VV | I | helix |
| 255 | = H HH | HH | A | active site |
| 256 | I LL | LL | s | helix |
| 257 | L ML | LN | I | helix |
| 258 | C SS | AA | I | helix |
| 259 | H HH | HH | s | helix |
| 260 | I KK | RF | S | helix |
| 261 | K KK | LF | I | helix |
| 262 | FG GG | AA | s | helix |
| 263 | AE AT | EE | I | helix |
| 264 | VA AS | IV | s | helix |
| 265 | SW VV | WW | s | helix |
| 266 | N HH | HH | . | helix |
| 267 | LAL LR | LL | i | helix |
| 268 | VP FF | LL | . | helix |
| 269 | KRR TR | RR | s | helix |
| 270 | K NM | TT | I | helix |
| 271 | K KM | VV | s | helix |
| 272 | | YY | . | helix |
| 273 | QQ | RR | . | |
| 274 | GA AK | AA | s | |
| 275 | NN SY | QQ | i | parse |
| 276 | LII PK | QQ | I | parse |
| 277 | DSS GG | GG | s | parse |
| 278 | VA VA | GG | s | parse |
| 279 | EQ IK | KK | s | |
| 280 | III II | II | S | parse |
| 281 | GGG GG | SS | I | parse |
| 282 | TTT VV | TT | I | parse |
| 283 | TA AA | VV | I | |
| 284 | NN RR | II | i | |
| 285 | LTV II | SS | . | |
| 286 | ME SH | TT | s | |
| 287 | TS SH | RR | . | parse? |
| 288 | PY WS | WW | s | |
| 289 | VH AA | FF | I | |
| 290 | YY VV | LL | s | parse? |
| 291 | LP PY | PP | . | |
| 292 | QA YY | RR | S | parse |
| 293 | TS RS | DD | s | parse |
| 294 | EE RT | CS | S | parse |
| 295 | R TS | DD | . | parse |
| 296 | | SS | I | |
| 297 | G Y | | . | |
| 298 | VH | KK | . | |
| 299 | KK | QQ | s | |
| 300 | VA KE | BB | I | parse |
| 301 | SE ED | PP | s | |
| 302 | ED DD | CS | I | |
| 303 | II MK | II | . | |
| 304 | EE EA | VV | s | |
| 305 | RA CA | AA | I | |
| 306 | AA AA | RK | . | |
| 307 | LB LA | EE | I | |
| 308 | VL RR | MM | s | |
| 309 | SS VT | YY | . | |
| 310 | NI P | NN | I | |
| 311 | P | VV | I | |
| 312 | KK | DD | s | strand |
| 313 | OS GS | QQ | s | strand |
| 314 | ML LF | FF | . | strand |
| 315 | DA SH | MM | I | |
| 316 | NG GS | GG | I | |
| 317 | QR DD | WW | I | |
| 318 | LW WW | AA | s | |
| 319 | FY YF | FF | I | parse |
| 320 | LL LF | AA | A | parse |
| 321 | ED PP | HH | I | |
| 322 | PQ PP | EE | I | |
| 323 | VV II | LL | I | |
| 324 | LL YY | TT | . | |
| 325 | FQ K | FF | s | |
| 326 | | KK | | |
| 327 | G GG GG | GG | | |
| 328 | K GR ES | RR DD | | |

Column 4 (residues 407–484):

```
407  TT HI SS            s   parse
408  DN TT ST            .   parse
409  SP AA FV            a
410  LM LM DD            s   parse
411  IT MM AA            i
412 D NN DD DD           s   parse
413  AI AA RR            S   parse
414  GG GG               -i
415  VV VV               S i
416  DK AA               s   parse
417  LL SS               I i
418  TS TT II            S   parse
419  FF FY AV            S   parse
420  EE ND DD            s
421  HK NN RR            -i
422  MH SS WW            s   parse
423  RR RR               s   parse
424  GG GG PP            i
425  KI EE SS            -i
426  PP YF SS            S   parse
427  LL EL GG            i   parse
428  AA FY LL            S   parse
429  EE EE KK            -i
430  MR VL WW            s
431  AA FF LL            s/i strand (bl)
432  SS EE MM            -i  strand (bl)
433                      s   strand (bl)
434                      i   strand (bl)
435                      s   strand (bl)
436                      s/i strand (bl)
437                      s
438                      s
439                      -i
440                      s
441                      s
442  ST                  i
443  WW K                -i
444  CY DN               s
445  IV AG               s
446  VV KV KV            s   helix
447  YY NN TT            i   helix
448  PP SS PF            s   helix
449  -M YY GG            s/i helix
450  FF YY FF            i   helix
451  KK RR               s   helix
452  EGG II              i   helix
453  DI LL LL            s   helix
454  YY NN               s   helix
455  VV LI               i   helix
456  FF VV               s   helix
457  MM KK               s   helix
458  CC DD EE            i
459  IIF FF YY           S   parse
460  LK TT               .   parse
461  IK KK               s   parse
462  I I                 i   strand
463  TN GG NN            a   strand
464  ND NN               i   strand
465  QL LL II            s   strand
466  PN TD EE            A   active site
467  VK PP               A   active site
468  SI YY               A   active site
469  TT YY               .
470  TS TT TT            i
471  EE EE EE            s
472  NN NN               d
473  GG GG GG            a
474  AA AA               s
475  IF VV               i
476  TT IM               s
```

(The page is a dense multi-column protein secondary-structure prediction listing; residue numbers run roughly 283–484 across four columns, each with amino-acid triplets and structure annotations such as "parse", "helix", "strand", "coil", "break to inside", "break to surface", "active site", "amphiphilic", "distributed parse".)

Fig. 1. Multiple alignment of the two subfamilies of the 6-phospho-strand-D-galactosidase. Undefined residues correspond to parses. In regions where the alignment has been readjusted by hand, surface and interior assignments may not correspond to those produced by the automated computer output. These should be ignored.

Subfamily a (a b c j k o):

a - (p11546) lacg_lacla 6-phospho-strand-galactosidase (EC 3.2.1.85) (strand-D-phosphogalactoside galactohydrolase). *Lactococcus lactis* (subsp. *lactis*) (*Streptococcus lactis*).

b - (p11175) lacg_staau 6-phospho-strand-galactosidase (EC 3.2.1.85) (strand-D-phosphogalactoside galactohydrolase). *Staphylococcus aureus*.

c - (p14696) lacg_lacca 6-phospho-strand-galactosidase (EC 3.2.1.85) (strand-D-phosphogalactoside galactohydrolase) (p-strand-gal) (pbg). *Lactobacillus casei*.

d - (p24240) ascb_ecoli 6-phospho-strand-glucosidase (EC 3.2.1.86). *Escherichia coli*.

e - (p26206) arbb_erwch 6-phospho-strand-glucosidase (EC 3.2.1.86). *Erwinia chrysanthemi*.

f - (p11988) bglb_ecoli 6-phospho-strand-glucosidase (EC 3.2.1.86). *Escherichia coli*.

Subfamily b (d e f g h i l m n p q r):

a - (p26208) bgla_clotm strand-glucosidase a (EC 3.2.1.21) (gentiobiase) (cellobiase) (strand-D-gluco-side glucohydrolase). *Clostridium thermocellum*.

b - (p10482) bgls_calsa strand-glucosidase (EC 3.2.1.21) (gentiobiase) (cellobiase) (strand-D-gluco-side glucohydrolase) (amygdalase). *Caldocellum saccharolyticum*.

c - (p22073) bgla_bacpo strand-glucosidase a (EC 3.2.1.21) (gentiobiase) (cellobiase) (strand-D-gluco-side glucohydrolase) (amygdalase). *Bacillus polymyxa*.

d - (q03306) bgla_bacci strand-glucosidase (EC 3.2.1.21) (gentiobiase) (cellobiase) (strand-D-glucoside glucohydrolase) (amygdalase). *Bacillus circulans*.

e - (p22505) bglb_bacpo strand-glucosidase b (EC 3.2.1.21) (gentiobiase) (cellobiase) (strand-D-gluco-side glucohydrolase) (amygdalase). *Bacillus polymyxa*.

f - (p12614) bgls_agrsp strand-glucosidase (EC 3.2.1.21) (gentiobiase) (cellobiase) (strand-D-glucoside glucohydrolase) (amygdalase). *Agrobacterium sp.* (strain atcc 21400).

g - (q00326) myro_brana myrosinase precursor (EC 3.2.3.1) (sinigrinase) (thioglucosidase). *Brassica napus* (rape).

h - (p09849) lph_rabit pos 1361 to 1926 of lactase-phlorizin hydrolase precursor (EC 3.2.1.108) (EC 3.2.1.62) (lactase-glycosylceramidase) (lph). *Oryctolagus cuniculus* (rabbit).

i - (p29092) myr3_sinal myrosinase mb3 precursor (EC 3.2.3.1) (sinigrinase) (thioglucosidase). *Sinapis alba* (white mustard).

j - (p26204) bgls_tripp noncyanogenic strand-glucosidase precursor (EC 3.2.1.21). *Trifolium repens* (creeping white clover).

k - (p09849) lph_human pos 1361 to 1927 of lactase-phlorizin hydrolase precursor (EC 3.2.1.108) (EC 3.2.1.62) (lactase-glycosylceramidase). *Homo sapiens* (human).

l - (p26205) bglt_tripp cyanogenic strand-glucosidase precursor (EC 3.2.1.21) (linamarase) (fragment). *Trifolium repens* (creeping white clover).

TABLE 1. Consensus Secondary Structure Prediction for the 6-Phospho-β-D-galactosidase Superfamily*

| | | | | |
|---|---|---|---|---|
| Strand A | 009–011 | Strand A | 049–051 | |
| Strand B† | 014–020 | Strand B† | 053–060 | Internal |
| Helix 1† | 072–084 | Helix 1† | 095–107 | Amphiphilic |
| Strand C? | 089–093 | Strand C | 111–115 | Amphiphilic |
| Act site a | 095–102 | Act site a | 117–125 | |
| Helix 2† | 116–130 | Helix 2† | 138–153 | Amphiphilic |
| Strand D† | 136–140 | Strand D† | 159–163 | Internal |
| Act sit b† | 141 | Act sit b† | 164–166 | |
| Helix 3† | 158–177 | Helix 3† | 181–198 | Amphiphilic |
| Strand E | 182–185 | Strand E | 205–208 | |
| Act sit c | 184–187 | Act sit c | 207–209 | |
| Helix 4† | 212–226 | Helix 4† | 248–268 | Largely internal |
| Strand x | 242–245 | | | Ambiguous alignment |
| Helix x | 259–273 | | | Shifted alignment |
| Strand y | 275–280 | | 318–320 | Amphiphilic |
| Helix 5† | 286–293 | Helix 5† | 332–342 | Interior |
| Strand F | 323–327 | Strand F | 381–388 | Ambiguous alignment |
| Helix y | 329–339 | Gap | | |
| Strand z | 375–382 | Strand z | 446–452‡ | Amphiphilic |
| Helix 6† | 385–398 | Helix 6† | 456–469 | Amphiphilic |
| Strand G | 404–407 | Strand G† | 476–479 | Internal |
| Act site d† | 408–410 | Act site d† | 480–482 | |
| Helix 7† | 431–448 | Helix 7† | 497–517 | Amphiphilic |
| Strand H† | 450–454 | Strand H† | 521–525 | Amphiphilic |
| Strand I† | 456–459 | Strand I† | 527–530 | Interior |
| Strand J† | 464–467 | Strand J† | 535–539 | Interior |
| Strand K† | 478–482 | Strand K† | 548–554 | Interior |
| Helix 8† | 496–509 | Helix 8† | 563–576 | Amphiphilic |

*Assignments in the consensus model (which applies to the entire superfamily) are designated with upper case letters A–K (for β strands) and Arabic numerals 1–8 (for α helices). Strands and helices designated by "x," "y," and "z" are not part of the consensus model, and may be present in only some members of the superfamily. Assignments marked with "?" are weak within one subfamily, but confirm a stronger assignment in the other subfamily.

†Reliable assignments.

‡The multiple alignment is ambiguous; see text.

Strand H (a450–454; b521–525) is amphiphilic and confirmed in both subfamilies.

Strand I (a456–459; b527–530) is interior and confirmed in both subfamilies. It may be longer by two residues in subfamily b.

Strand J (a464–467; b535–539) is interior and confirmed in both subfamilies.

Strand K (a478–482; b548–554) is interior, well anchored, and confirmed in both subfamilies.

Finally, helix 8 (a496–509; b563–576) is amphiphilic, well anchored, and confirmed in both subfamilies.

In examining the consensus secondary structural model reported in Table 1, it is difficult not to notice the secondary structural pattern characteristic of an 8-fold α–β barrel protein. This tertiary structural hypothesis does not rest solely on pattern recognition. The model is, in fact, enforced by the active site assignments designated in Table 1. Here, β strands C, D, E, and G all must terminate near the active site of the protein, as in an 8-fold α–β barrel. While other topologies could also bring these residues together, this is our preferred tertiary structural model.

## REFERENCES

1. Benner, S.A. Patterns of divergence in homologous proteins as indicators of tertiary and quaternary structure. Adv. Enzyme Reg. 31:121–181, 1989.
2. Pananoyotou, G., Bax, B., Gout, I., Federwisch, M., Wroblowski, B., Dhand, R., Fry, M.J., Blundell, T.L., Wollmer, A., Waterfield, M.D. Interaction of the p85 subunit of PI 3-kinase and its N-terminal SH2 domain with a PDGF receptor phosphorylation site. EMBO J. 11:4261–4272, 1992.
3. Russell, R.B., Breed, J., Barton, G.J. Conservation analysis and structure prediction of the SH2 family of phosphotyrosine binding domains. FEBS Lett. 304:15–20, 1992.
4. Musacchio, A., Gibson, T., Lehto, V.-P. and Saraste, M. SH3—an abundant protein domain in search of a function. FEBS Lett. 307:55–61, 1992.
5. Bazan, J.F. Structural design and molecular evolution of a cytokine receptor superfamily. Proc. Natl. Acad. Sci. U.S.A. 87:6934–6937, 1990.
6. Moe, G.R., Koshland, D.E., Jr. In: "Microbial Energy Transduction, Genetics, Structure and Function of Membrane Proteins." Youvan, D.C., Daldal, F., eds. Cold Spring Harbor, NY: Cold Spring Harbor Press, 1986: 163–168.
7. Benner, S.A., Gerloff, D.L., Jenny, T.F. Predicting protein crystal structures. Science 265:1642–1644, 1994.
8. Benner, S.A., Ellington, A.D. Evolution and structural theory: The frontier between chemistry and biochemistry. Bioorg. Chem. Frontiers 1:1–70, 1990.
9. Chothia, C., Lesk, A.M. The relation between the divergence of sequence and structure in proteins. EMBO J. 5:823–826, 1986.
10. Benner, S.A., Gerloff, D.L., Jenny, T.F., Chelvanayagam,

G., Knecht, L.J., Cohen, M.A., Gonnet, G.H. Predicting conformation in proteins from homologous protein sequences. Nature Struct. Biol., submitted.

11. Jenny, T.F., Benner, S.A. Evaluating predictions of secondary structure in proteins. Biochem. Biophys. Res. Commun. 200:149–155, 1994.

12. Jenny, T.F., Benner, S.A. A prediction of the secondary structure of the pleckstrin homology domain. Proteins: Struct. Funct. Genet. 20:1–4, 1994.

13. Musacchio, A., Gibson, T., Rice, P., Thompson, J., Saraste, M. The pH domain—a common piece in the structural patchwork of signaling proteins. Trends Biochem. Sci. 18: 343–348, 1993.

14. Gerloff, D.L., Jenny, T.F., Knecht, L.J., Benner, S.A. A secondary structure prediction of the hemorrhagic metalloprotease family. Biochem. Biophys. Res. Commun. 194: 560–565, 1993.

15. Bairoch, A., Boeckmann, B. The SWISS-PROT protein sequence data bank. Nucleic Acids Res. 20:2019–2022, 1992.

16. Gonnet, G.H., Benner, S.A. Computational Biochemistry Research at ETH. Technical Report 154, Departement Informatik, March, 1991.

17. Gonnet, G.H., Cohen, M.A., Benner, S.A. Exhaustive matching of the entire protein sequence database. Science 256:1443–1445, 1992.

18. Benner, S.A., Badcoe, I., Cohen, M.A., Gerloff, D.L. Bona fide prediction of aspects of protein conformation: Assigning interior and surface residues from patterns of variation and conservation in homologous protein sequences. J. Mol. Biol. 235:926–958, 1994.

19. Benner, S.A., Gerloff, D. Patterns of divergence in homologous proteins as indicators of secondary and tertiary structure: The catalytic domain of protein kinases. Adv. Enzyme Reg. 31:121–181, 1991.

20. Wacker, H., Keller, P., Falchetto, R., Legler, G., Semenza, G. Location of the two catalytic sites in intestinal lactase-phlorizin hydrolase. J. Biol. Chem. 267:18744–18752, 1992.