The Phospho-β-galactosidase and Synaptotagmin Predictions

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

ABSTRACT Two bona fide consensus predictions of secondary and tertiary structure in a protein family, made and announced before experimental structures were known, are evaluated in light of the subsequently determined experimental structures. The first, for phosphoβ-galactosidase, identified the core strands of an 8-fold α - β barrel, and identified the 8-fold α - β barrel itself, which was found in the subsequently determined experimental structure to be the core folding domain. The second, for synaptotagmin, identified seven out of eight β-strands in the structure correctly, missing only a noncore strand. Three preferred "topologies" were selected from several hundred thousand possible topologies of these seven predicted strands using a rule-based analysis. The subsequently determined experimental structure showed that these seven strands in synaptotagmin adopt one of the three preferred topologies. We were unable, however, to identify the correct topology from among these three topologies. © 1995 Wiley-Liss, Inc.

INTRODUCTION: SELECTION OF STRUCTURES

The ETH group prepared three predictions for the contest consummated in Asilomar. These were done in two phases. First, predictions for phospho-β-galactosidase^{1,2} and isopenicillin N-synthase³ were prepared. These proteins were chosen because they were large, mechanistically interesting enzymes with sufficient members in the evolutionary family to permit a consensus model to be built. In September, we were advised that the structure for isopenicillin N-synthase would not be solved in time for the contest, and asked to submit another prediction. At that time, only synaptotagmin (or C2 homology domain protein)^{4,5} remained as a suitable target. Therefore, we prepared a prediction for synaptotagmin as well.⁶ The synaptotagmin and phospho-β-galactosidase predictions are discussed below. An experimental structure for isopenicillin N-synthase is still not available.

SUMMARY OF THE PREDICTIONS

For both synaptotagmin and phospho- β -galactosidase, a residue-by-residue secondary structural model was first predicted.^{2,6} Tertiary structural models were then built by assembling the predicted secondary structural elements. With phospho- β -galactosidase, the predicted secondary structure elements,² oriented using assignments of active site residues, suggested that the predicted secondary structural units assemble to form an 8-fold α - β barrel as the core domain.²

With synaptotagmin,⁶ a combinatorial analysis was used to build a set of all possible β -sheet packing orientations for the seven B-strands that were predicted in the structure.7 Following semiempirical rules, the number of possible folding "topologies" was reduced from several hundred thousand to just three. These corresponded to the β -sheet topologies found in the retinol binding protein, pseudoazurin, and the pleckstrin homology domain. Coordinate representations (see Figures 4, 5, and 6 in reference 6) of these three possible folds (or "topologies") were built based on the crystal structures of these domains or, in the case of the pleckstrin homology domain, on a model built from a published description of the structure.⁶ Further, more speculative rules were then use to suggest that the pleckstrin homology domain model was best among the three preferred topologies.

EVALUATING THE PREDICTIONS

As a rule, we do not evaluate our own predictions until after they have been evaluated by others. Therefore, we are indebted to Thomas Defay and Fred Cohen for having served as judges in this contest.⁸ Their insightful evaluation, appearing elsewhere in this issue, has enabled us to focus on a few points, important to a "right, wrong, why, learned" discussion of structure prediction generally, but illustrated well by these two predictions. Figures 1 (synaptotagmin) and Figure 2 (phospho- β -galactosidase) summarize the comparison between secondary structure as predicted and as assigned to crystallographic data.⁹ These figures serve as reference for the discussion below.

Received April 18, 1995; revision accepted August 21, 1995. Address reprint requests to Steven A. Benner, Department of Chemistry, ETH Zürich, CH-8092 Zürich, Switzerland.

| | • | - | | | | | | | | | | | | | | | - | | | |
|-----|----------|--------|--------------|------------|-----------|--------------------|------------------|---|--------------|--|------------|--------------|----|--------|-------------|------------|------------|-------|-------------|------|
| | _ | | | - | | 0.000 000 0 | | | 1. | fadaba | ~ | v | Ŧ | -11 | Brock | ti at ad | Evnerin | ontal | Cat+ | |
| | Pos | 5 q | ΙW | гуля | csurv | DCBEFGH | р | nomikji | n | reacba | g I | Г | Ŧ | 0 | rret Car | | Experim | | binding | je i |
| | | | | | | | | | | | | | | | | surf | | | prnurnd | |
| | | | | | : | | | | | | | | | • | ETH | inter | access | Str | - | |
| | 1 | ιc | C | CCCC | CCCCC | CCCCCCC | | 1. A. | | | · 1 | | | | | | | | | |
| | 1 2 | 2 G | G | GGGG | GGGGG | GGGGGGG | S | GKEEEEE | Н | OKKKKK | ΚI | Р | F | F | | S | | | · _ · | · |
| | | ۲ ۲ | í ř | VVVV | MMMMM | TTTTTTTT | E | GEEEEE | 0 | SEEEEE | ΕI | т | N | N | | S | 1 A. | | | • |
| | 7 | | | חחחח | ממממת | DDDDDDD | 2 | 0000000 | ř | FFPDFF | F I | s | S | S | | S | | | | |
| | 4 | | , | סטעעע | | ן עעטעעט | . г. В | QAFFQQQ | 17 | OFFEFE | | | | | | s | 185 | | | |
| | | | . н | нннн | кннин | нннннн | | LEEEEE | | QEEEEE | V 1 | 5 | 7 | 7 | | | | | | |
| | | | | | | TTTTTTT | | | | | | | | | | s | 192 | | | |
| | | | | | | EEEEEE | L | | С | LLLLLL | L (| Ŕ | | | • | · · | 32 | | · . | |
| | 8 | 3 V | 7 R | RRRR | KKKKK | RRRRRRR | | Ľ. | | | I | Ι | т | Α. | | | • | | - | |
| | ģ | R | R R | RRRR | RRRRR | RRRRRRR | _ | | | · · · · | I | ĸ | 0 | H | | · · · · | | | | |
| | - | | | | | | ē | GGGGGGG | G | GGGGGG | Γi | S | Ĝ | G | | | 6 | | | |
| | | | | | | RRRRRRR | | | | | | | P | | в | S | 69 | в1 | | |
| | | | | | | IIIIIII | | | | | | | ŵ | | Ē | ī | Õ | B1 | | |
| | | | | | | | | | | | | | | | | | .19 | BI | | |
| | | | | | | YYYYYYY | | | | | | _ | W | -44 | В | S | | | | |
| | | | | | | | | FFTTFFF | | | | | _ | _ ' | В | i | 3 | B1 | | |
| | | | | | | 0000000 | | SSSSSSS | | | | | _ | _ | в | S | | B1 | | |
| | 16 | V V | ' I | IIII | ААААА | AAAAAAA | \mathbf{L} | LLLLLLL | L | LLLLLL | L | · | _ | _ | в | I | .2 | B1 | | |
| | 17 | E | N | RRRR | EEEEE | нннннн | C, C | RRRRRRR | R | EDDDDD | D [| | | | в | S | 45 | B1 | * | |
| | 18 | L | v | AAAA | VVVVV | IIIIII | Y | YYYYYYY | Y | YYYYYY | Y | _ | _ | _ | в | I | 4 | B1 | | |
| | 19 | ĸ | к | PPPP | ATTTT | EDDDDDD | \mathbf{L} | vvvvvvv | т | DDDDDD | DI | | | | | i | 29 | B1 | * | |
| | 20 | | | | | | | PPPPPPP | | | | - | | | , | | 93 | | | |
| | 21 | _ | · | TTTT | | | | TTTTTTT | | | | х | R | m | | s | 164 | | | |
| | | | = | | DDDDD | RRRRRR | | | | | | | P | | | S | 104 | | | |
| | | | | | | | | AAAAAAA | | | | | | | • | | | | | |
| | | | | | | EEEEDD | | GGGGGGG | | | | | Е | | _ | S | 85 | | | |
| | 24 | Ν | \mathbf{L} | EEEE | KKKKK | vvvvvvv | | KKKKKKK | _ | | _ | R | Ŗ | R | В | S | 36 | B2 | | |
| | 25 | L | Ŀ | IIII | LLLLL | LLLLLLL | Ŀ | LLLLLLL | Ĺ | LLLLLL | L | \mathbf{L} | Ŀ | L | в | I · | 0 | в2 | | |
| | 26 | Κ | т | нннн | ннннн | IIIIIII | т | TTTTTTT | v | ATTILL | Т | I | R | N | в | S | 36 | B2 | | |
| | | | | | | vvvvvvv | т | | v | $\nabla \nabla \nabla \nabla \nabla \nabla \nabla$ | v | v | v | v | в | I | 0 | B2 | | |
| | | | | | | VVVVVLL | | VVCCVVV | | | | | R | | B | S | . 8 | B2 | | |
| | | | - | | | | | | | | | | I | | B | I | 0 | B2 | | |
| | | | | | | ~~~~ | | IIIIIII | | | | | | | | T | - | | | |
| | | | | | | RRRRRRR | | LLLLLLL | | | | | I | | В | | 33 | B2 | | |
| | 31 | Е | E | EEEE | DDDDD | DDDDDDD | К | EEEEEE | К | 000000 | Q | | S. | | b | | 61 | B2 | | |
| | | | | | | АААААА | А | ААААААА | А | аааааа | А | | G | | b | i | 1 | B2 | | |
| | -33 | Α | R | RRRR | KKKKK | KKKKKKK | \mathbf{T}^{2} | KKKKKKK | \mathbf{r} | EAAAAA | E | R | Q | Q | b | S | 16 | в2 | | |
| ۰., | 34 | N | N | NNNN | NNNNN | NNNNNN | Ν | NNNNNN | D | EEEEEE | D | Q | Q | Q | | S | 117 | | | |
| | 35 | L | L | LLLL | LLLLL | LLLLLL | L | LLLLLL | L | LLLLLL | L | L | L | L | | I | - 7 | | | |
| | 36 | Ι | Ŧ | IIII | IIIII | vvvvvvv | к | KKKKKKK | Р | PPPPPP | Ρ. | ·P | P. | Ρ | | • | 36 | | | |
| | | | | | | PPPPPPP | | KKKKKKK | | | | | ĸ | | | | 47 | | - | |
| | | | | | | MMMMMM | | MMMMMM | | | | | | v . | | I | 77 | | | |
| | | | | | | | | | | | | | | | | <u> </u> | | | | , |
| а | | | | | | DDDDDDD | _ | DDDDDDD | _ | | | | N | | | 5 | 49 | | | |
| | | | | | | PPPPPPP | _ | VVVVVVV | | | | | K | | | | 192 | | | |
| | • 41 | Ν | Ν | NNNN | NNNNN | NNNNNN | т | GGGGGGG | N | GGGGGG | S | S | Ν | N | 1. . | S | 66 | | | |
| | 42 | | | | | | | | | | | Т | | | . 5 | | | | | |
| | 43 | _ | | | | | _ | | _ | | _ | к | K | K | | | | | | |
| | 44 | | | | <u></u> | | _ | | | · · | - . | | N | | - e | | | | | |
| | 45 | 5 | ā | 0000 | 22222 | GGGGGGG | ā | 6666666 | <u>त</u> | 000000 | Ā | | s | | | | 44 | | | |
| | | | | | | LLLLLLL | | | | | | | ĩ | | | I | 29 | • | | |
| _ | | | | | | | | | | | | | | | | | | | | |
| а | | | | | | SSSSSSS | | | | | | | v | | | i | 0 | | · · · · | • • |
| | | | | | | | D | DDDDDDD' | D | DDDDDD | D | | D | | | · A . | 21 | | Ca++ | |
| | 49 | Ρ | P | PPPP | PPPPP | PPPPPPP | P٠ | PPPPPPP | P | PPPPPP | | | Р | | | - | 0 | B3 ' | | |
| | 50 | Y | Y | YYYY | YYYYY | YYYYYY | Υ· | YYYYYYY | Y | YYYYYY | Y | Ϋ́Υ | к | К | в | | 6 | в3 | ' | |
| | 51 | Ι | v | VVVV | VVVVV | VVVVVVV | v | VVVVVVV | v | VVVVVV | v | v | v | v | в | I | 0 46 | в3 | | |
| | | | | | | | | KKKKKKK | | | | | Ι | | В | s | 46 | B3 | | |
| | | | | | | | | IIIIIII | | | | | v | | В | I | 0 | 23 | | · · |
| | | | | | | | | | | | | | | | | | | | | |
| | | | | | | KKKKKKK | | | | | | | E | | B | | 48 | | | |
| | | | | | | LLLLLLL | | | | | | I· | | I | B B | Ι. | 0 | | | |
| | | | | | | IIIIIII | | MMLMMM | | | | v | Н | н | в | | | в3 | ÷ | |
| | 57: | P | Ρ | PPPP | PPPPP | PPPPPPP | | | | | | G | G | G | | | 104 | | 1. . | |
| | 58 | D | D | DDDD | DDDDD | DDDDDDD | D | NGNNNNN | D | DDDDDE | E : | т | V | v | | S | 104 74 | | | |
| | | | | | | | | GGGGGGG | | | | | | | | S S | 32 | | | |
| | 60 | S | ĸ | סמממ | KKKKK | KKKKKKK | P | KKKKKKKK | - • | | ••• | F | | | | S | 22 | | | |
| | 60 61 | 2 | D | AVANALAR . | NINTATATA | CCCCCCC | ц. 1 | TODDDDDD | | | - | | | | | | | | · · · · | |
| | | | | | | | | RRRRRRR | | | | D | | | | S | - / | | | |
| | | | | | | EEEEEE | | | | | | Q | | | | t = -t | 94 | | | |
| | | | | | | | | KKKKKKK | | | | ĸ | | | | | 38 | | | |
| | | | | | | KKKKKKK | К | KKKKKKK | К | KKKKKK | ĸ | v. | s | S | | | 77 | | | |
| | 65 | Κ | Κ | QQQQ | QQQQQ | 0000000 | R | KKKKKKK | F٠ | FYYYFF | v | Е | | | | S | 56 | B4 | | |
| | 66 | Κ | К | KKKK | KKKKK | KKKKKKK | | | | | | ĸ | | | | s | | B4 | | |
| | 67 | Т | т | TTTT | TTTTT | | | TTTTTTT | | | | Т | | | h | | 8 | | | |
| | 68 | к | R | RRRK | ккккк | KKKKKKK | s | SSTTTTT | ĸ | KKKKKK | ĸ | ĸ | | _ م | b b | | 154 | | <i></i> | |
| | 69 | т | Ť | TTTT | ጥጥጥጥ | TTTTTTT | T | VIVVITT | v | tanana. | | | | и. | b | а. Т | | | | |
| | | ÷ | - | | | | т | ATAATT | Y | V V V V V | v | v | v | V · | a | . Т | ,39 | | | |
| | | | | · | | • | | | | | | | | | | | | | | |

Ş.

.....

Fig. 1. Residue-by-residue secondary structure prediction for synaptotagmin. From left to right, the columns are alignment position number, multiple alignment using the one-letter code for amino acids, the predicted secondary structure, the predicted surface and interior assignments, the experimental surface and interior assignments, the experimental assignment of secondary structure, and the residues involved in calcium binding (asterisks refer to the "degenerate" calcium binding site.⁹) See Gertoff et al.⁶ for details. (Continued on overleaf.)

| 70IIVVVVIIIIIIIIIIIIKKKKKKKKHHQHHHHH71QKKKKKRRRRRKKKKKKKKKKRRR< | Υ Τ ΝΝΟGFΝΡWWD FNPPWWD FEEFE FEEFE FEV A Y | B B B B B B B | I S I S | 53 114 133 38 41 95 4 75 12 121 39 122 21 48 9 180 11 54 | 85 85 85 85 85 85 85 85 85 85 85 | · · · · · · · · · · · · · · · · · · · |
|--|--|--|------------------------------|--|--|---------------------------------------|
| 90 P P PPPP PPPP EEEEEEE N FFFFFFF F YYYYY F 91 Q E GGGG SSSSS SSSSSS E EEEEEEEE N AQQSSS N 92 D D DDDD DDDDD DDDDDDD N QQQQQQQ E DEEEEE E 93 R K VVVV KKKKK KKKKKK M MIIIIII L ALLLLL I 94 D D EEEE DDDDD DDDDDDD E QQQQQQQ K TTTTTT T 96 K R RRRR RRRR RRRRRR N CSQQQQQQ K TTTTTT T 97 R R RRRR RRRRR RRRRRR N CSQQQQQQ K TTTTTT T 98 L I LLLL LLLL LLLLL V LLVVVVV L LLLLLL L 99 L L SSSS SSSSS SSSSSS I VMCVVVV V LLLLLLL L 99 L L SSSS SSSSS SSSSSS I VMCVVVV V VVVVV V 100 I I VVVV VVVVE VVVVVV I VIVVVVV V VIVVVVV V 101 E E EEEE EEEEE EEEEEE A TTTTTT S AAAAAA A 102 V V VVVV IIIIII IIIIII V VVVVVV V IVVVVV V 103 W WWWW WWWW WWWWWW M VMLLLLL Y FYYYY Y 104 D D DDDD DDDDD DDDDDDD D 105 W W WWWW WWWWW WWWWWW Y YYYYYY F FFFFFF F 106 D D DDDD DDDDD DDDDDDD D 107 R R RRRR RRRRR LLLLLLL C RKKKKK R RRRRRR R 108 T T TTTT TTTTTTT I ILLILII F FFFFFF F 109 S S SSSS TTTTT SSSSSSS G GGGGGGG S SSSSS S 110 R R RRRR RRRRR RRRRRR H TSKKKKK R KKKKKK K 111 N N NNNN NNNNN NNNNNN N SNNNNNN H HHHHHH H 122 D D DDDD DDDDD DDDDDDD E EDDEDDD D DDDDDD | M L V L R F F V V V R F F V V V D E E D D D Y D D K S S G S S K K N N R D D | B B B B B B b b B B | sss ssssi si A Asi sssii sis | 35 0 8 0 4 0 48 4 33 51 161 187 68 155 52 80 80 13 2 46 2 82 6 | hhhh BBBBBBBB BBBBBB BBBBBBB BBBBBBBBBB | Ca++ Ca++ |
| 122 | N S S I L L R K Q Q G G G - Y Y R R H I H H I L K L L S K N N K N F N N | | 5 | 112 H 68 H 102 H 68 H | 38 38 38 38 38 38 | * |

Fig. 1. (Continued from overleaf.)

PHOSPHO-β-GALACTOSIDASE AND SYNAPTOTAGMIN

| | MTKTLPKDFIF <u>GGAT</u> AAYQAEGATHTDGKGPVAWDKYLEDNYWYTAEPAS | 50 |
|---|--|------------|
| | eeeeeeee hhhhhhh | exp |
| | eee eeeee hh | pred |
| | core strand 1 | |
| | | 100 |
| | DFYHKYPVDLELAEEYGVNGIRISIAWSRIFPTGYGEVNEKGVEFYHKLF | 100 |
| | hhhhhhhhhhhh eeeeeeeee hhhhhhhhh hhhhhhhh | exp |
| | core helix A core strand 2 core helix I | pred |
| | core nellx A core strand 2 core nellx h | . . |
| | AECHKRHVEPFVTLHHFDTPEALHSNGDFLNRENIEHFIDYAAFCFEEFP | 150 |
| | hhhhhh eeeee hhhhhh hhhhhhhhhhhhhhhh | exp |
| | hhhhh eeeee hhhhhhhhhhhhhhhhhhhhhhhhhh | |
| | core strand 3 core helix C | • . |
| ; | | |
| | EV <u>NYWTTF</u> NEIGPIGDGQYLVGKFPPGIKY <u>DLAKVFOSHHNMMVSHARAV</u> | 200 |
| | eeeee hhhhhhhhhh h eeee hhhhhhhhhhhhhhh | |
| | h eeee hhhhhhhhhhhhhhhhhhh | pred |
| | core strand 4 core helix D | |
| • | | |
| | <u>KLYKDKG</u> YK <u>GEIGVVHAL</u> PTKYPYDPENPADVRAAELFDIIHNKFILDAT | 250 |
| | hhhhh eeeeeeeeee hhhhhhhhhhhhhhhhhhhhh | exp |
| | hhhhh eeeeeeeeee hhhhhhhhhhhhhhhhhhhhh | pred |
| | core strand 5 core helix E | |
| ÷ | | 200 |
| | YLGHYSDKTMEGVNHILAENGGELDL <u>RDEDFOALDAAK</u> DLND <u>FLGINYY</u> M | |
| | hhhhhhhhhhhheeeeeee hhhhhhhhh eeeeeeee h | |
| | hhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhh | pred |
| | core stand o | |
| | SDWMQAFDGETEIIHSKYQIKGVGRRVAPDYVPWIIY | 350 |
| | | exp |
| | eeeeee ee eeee eeee eeee eeee | · • |
| | nnnnnnnn | prou |
| | PEGLYDOIMRVKNDYPNYKKIYITENGLGYKDEFVDNTVYDDGRIDYVKO | 400 |
| | hhhhhhhhhh eee eee hhhhhhhhh | |
| | hhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhh | |
| | core helix F core strand 7 core | F |
| | COTE METIX I COTE DEtana , | |
| | <u>HLEVLSDAIADGANVKGYFIW</u> SLMDVFSWSNGYEKRYGLFYVDFDTQERY | 450 |
| | hhhhhhhhh eeeeeeee eee eeee eeee | exp |
| | hhhhhhhhh eeee eee eee | pred |
| | helix G core strand 8 | - |
| | HELLY G COLC DELANG C | |
| | PKKSAHWYKKLAETOVIE | |
| | e hhhhhhhhhhh exp | • |
| | | |

e hhhhhhhhhhh exp hhhhhhhhhhhh pred core helix H

Fig. 2. Experimental structure assigned by DSSP: h indicates an α -helix, e indicates a β -strand. The underlined regions designate the core secondary structural elements in the conserved α - β barrel domain. These are assigned using the DEFINE program; the differences between the DSSP and DEFINE assignments (made on the same experimental data) are worthy of note. This figure illustrates the accuracy of the consensus prediction in the assignment of secondary structure to elements of secondary structure that are conserved throughout the protein family, but not (by definition) to those that are not. See Gerloff and Benuer² for details.

USING TERTIARY STRUCTURAL INFORMATION (SURFACE AND INTERIOR ASSIGNMENTS) TO ASSIGN SECONDARY STRUCTURE: ARE THE SECONDARY STRUCTURE PREDICTIONS CORRECT FOR THE CORRECT REASONS?

The protein structure prediction problem embodies a "chicken-or-egg" paradox that arises because tertiary structural interactions are stronger than local sequence interactions in determining secondary structure.¹⁰ This implies that predicting secondary structure from primary structure is essentially impossible without having at least some tertiary structural information.¹¹ At the same time, a reliable model for secondary structure appears to be necessary before predicting tertiary structure. In these predictions,^{12,13} this paradox is resolved by extracting tertiary structural information from a set of aligned homologous sequences before assigning a secondary structure.¹⁴ This tertiary structural information comes primarily in the form of surface and interior assignments, although active site assignments are also obtained (see below). Patterns in these assignments are then used as indicators of secondary structure.

This approach has had reasonable success in many laboratories in making bona fide predictions of secondary structure,¹⁵ those announced before an experimental structure is known. This success was mirrored in these predictions. For example, the synaptotagmin prediction identifies the first seven β -strands of the fold essentially correctly (Fig. 1). Further, with the exception of β 4, there is excellent correspondence between the beginnings and ends of the strands as predicted and as experimentally assigned by DSSP. These include all of the strands that Sutton et al. identify as the "core" strands.⁹

It is important to ask whether the close correspondence between the predicted and experimental assignments of secondary structure in the synaptotagmin family was achieved for the correct reasons. Figure 1 shows that it was. Both the predicted assignment of secondary structure (S and s for strong and weak surface assignments, I and i for strong and weak interior assignment) and the experimental assignments (from DSSP) are given. For $\beta 1$, $\beta 2$, $\beta 3$, $\beta 5$, and β 7, the correct assignment of secondary structure transparently arises from an accurate prediction of tertiary structural details, with the 2-residue periodicity of internal and external assignment characteristic of a B-strand.¹⁶ Thus, the physicochemical basis for the secondary structure assignment in these segments is clear. B4 is too short to be analyzed in this fashion with statistical significance. The segment containing $\beta 6$ was correctly identified as being largely internal, and the secondary structure correctly assigned using a different rule-based approach.

TRANSPARENCY IS USEFUL FOR IMPROVING A METHOD

There has been much discussion as to whether humans should be involved when predicting secondary structure, or whether a fully automated computer program should be the only goal of efforts in the field.¹⁷ At one level, this discussion is purely technological: Does the man (woman) or the machine produce the better results? Defay and Cohen have addressed the issue in this context.⁸

We frequently pointed out that the issue can also be addressed in terms of research strategy.^{11,12} Assume that we do not yet have a method that reliably assigns conformational information to a protein. Is it better at this point to design fully automated methods? Or to develop methods that involve human intervention? We have argued at some length that until the conformational problem is solved, human involvement in the bona fide prediction process is critical. Only through human involvement can one understand why a method works (when it works) and why it fails when it fails. It is this understanding that underlies rational improvement of prediction methods. Thus, a transparent method that performs poorly can often be more useful for developing the underlying science than a nontransparent method, even if it performs better.

For example, the synaptotagmin prediction is marred by a serious error: the final strand was misassigned as a helix (Fig. 1). With a transparent prediction tool, we are able to ask why. Inspection of the multiple alignment in Figure 1 shows that two factors figured significantly in the misassignment. First, a manual readjustment of the multiple alignment shifted a gap from the middle of the region mispredicted as a helix. As a gap is regarded as a "parsing" element, this shifting was necessary to predict the helix; had the gap not been shifted, a helix would not have been assigned in this region. Here, human intervention diminished the technological quality of the prediction. However, the transparency of the process allows us to understand what happened that led to the misassignment. In particular, the fact that the alignment is difficult to construct in this region suggests that secondary structure may have diverged in this region. This leads to further analysis below.

Second, the misassignment was aggravated because of an incorrect assignment of a conserved tryptophan at position 135 to the inside. This Trp residue lies on the surface in the experimental structure. This misassignment obscured a tertiary structural assignment pattern that would have identified this segment as a strand. It is interesting to note that in the broader multiple alignment provided by Sutton et al.,⁹ this Trp is not conserved, but rather displays a pattern of variation and conservation characteristic of a surface residue (but see below concerning the quality of the multiple alignment). A similar problem was encountered in the prediction of the SH3 domain,¹⁸ and future heuristics for assigning surface/interior residues will flag conserved tryptophan residues.

WHAT DOES IT MEAN WHEN ONE SAYS THAT "WHEN THE QUALITY OF AN ALIGNMENT IS POOR, PREDICTION ACCURACY SUFFERS"?

It is frequently noted that predictions based on multiple alignments suffer when the multiple alignment is poor. This comment creates the impression that if only we could improve tools for creating multiple alignments, than we would be better able to predict secondary structure. This impression is mis-

450

taken, in a way that is important for the future development of prediction tools.

A structural prediction made for a family of aligned homologous protein sequences is a consensus prediction. It is based on the assumption that proteins related by common ancestry all have similar conformation. This assumption is true only as an approximation. As proteins within a family diverge at the level of sequence, they also diverge at the level of conformation, first in loops, then in side chain orientations, then in secondary structure, and finally through the modification of conformation of entire domains. An excellent quantitative statement of this divergence was given by Chothia and Lesk.¹⁹

The breakdown of the assumption underlying all evolution-based structure prediction methods has an impact on the way in which predictions are evaluated.^{20,21} A consensus model can be accurate for only those structural features that a protein family holds in common.²¹ Further, as a consensus model is generally evaluated by examining only a single member of the protein family that it covers, at least some of the "mistakes" reflect not errors introduced by the prediction method, but rather cases where the individual experimental structure has a conformation different from the conformation of other proteins in the family. This represents, of course, an intrinsic limitation of the consensus modeling strategy.

Regions where the conformation has diverged within a family are, as a rule, regions that are difficult to align. Conversely, regions that are difficult to align are, quite often, regions where conformation has diverged. Regions that are difficult to align are therefore regions where the assumption underlying all evolution-based structure prediction tools is likely not to apply. Thus, these are regions where a consensus prediction is likely to be judged to be "wrong," especially when a single experimental structure is used to evaluate the prediction, not because the tool used to make the prediction was flawed, but because the strategy itself is based on an assumption that is not universal.

Fortunately, secondary structural elements that are not conserved rarely lie at the core of the folded structure. Thus, as noted by Defay and Cohen,⁸ the nonuniversality of the assumption does remarkably little damage when building tertiary structural models from a predicted consensus secondary structural model. For example, the phospho- β -galactosidase family contains two major subfamilies with low sequence similarity overall.² Large segments of the protein are difficult to align, implying that in these regions, secondary structure has diverged substantially. As a result, the single known experimental structure, obtained for protein "a" (Subfamily a) in the alignment,² is insufficient to evaluate the secondary structure in regions where secondary structure has widely diverged (Fig. 2); several of the elements predicted in the consensus model are simply wrong when evaluated using this particular experimental structure.⁸ Nevertheless, because core structural regions are conserved in both subfamilies of the phospho- β -galactosidase family and because these are the important elements when assembling a tertiary structural model, the consensus prediction (Fig. 2) is adequate as a starting point to build an essentially correct tertiary structural model for the "core" fold, the fold of the protein that both subfamilies adopt in common.

Likewise, the misassigned helix in the synaptotagmin family lies in a region where the alignment is extremely poor. This is, of course, the reason why the gap was misplaced in this region, and one reason why the secondary structural element was misassigned (see above). However, the information that we presently have does not allow us to conclude that a strand should be predicted in this region in the consensus model, as secondary structure within this region may not be conserved. We must wait for more experimental structures from the synaptotagmin superfamily before we will know whether this misassignment reflects a weakness in the secondary structure prediction heuristic, or whether it reflects divergence in secondary structure within the family, and therefore a weakness of the consensus modeling approach generally. Again, however, the misassignment of a noncore secondary structural element did not obstruct the assembly of the correct tertiary structure for synaptotagmin as one of three alternative packings (see below).

TERTIARY STRUCTURAL MODELING VERSUS THREADING

Given reliable tools for assigning secondary structure, the task remaining is to build a tertiary structural model. One approach is to attempt to "thread" a predicted secondary structural model on to proteins known in the database.²² Alternatively, one might attempt to assemble the predicted secondary structural elements de novo.

In the phospho-B-galactosidase prediction, the secondary structure pattern was obviously a signature of an 8-fold $\alpha-\beta$ barrel.²³ Further, the placement of active site residues in the model confirmed the β -barrel as the preferred topology for the overall fold. As the tertiary fold was identified in the context of a large number of known proteins with similar folds, the tertiary structure can be said to have been predicted by a "threading" approach in combination with a de novo secondary structure prediction. In this respect, this prediction is not distinctly different from the identification of a similar fold by a similar method by Kirschner and co-workers 7 years ago using a consensus secondary structure prediction obtained by the GOR method.^{24,25} As the reliability of secondary structure prediction tools

based on multiple sequence alignments has increased substantially since then, this sort of pattern recognition exercise will be routine in the future.

For synaptotagmin, the prediction of tertiary structure was more complicated. The secondary structure prediction did not in itself indicate a particular topology, only that the core structure was an all- β structure. There are, of course, many classes of β -sheets, barrels, and sandwiches, and neither the secondary structure prediction nor the active site assignments identified one of these as a clear template for building a tertiary prediction. Therefore, tertiary structural models were assembled de novo.⁶ First, a combinatorial approach was used to assemble all 322,560 possible sheet structures from the predicted β -strands. Then, a large majority of these were excluded by enforcing connectivity of strands in a β-sheet, avoiding loop crossovers, and using other rules that have (at least some) empirical basis.²⁶ This process reduced the number of possible β -sheet topologies first to 36, and then to six, grouped as three pairs of alternative folds. The details of this analysis are found in Gerloff et al.⁶

The database was then examined for analogs for the three folds remaining. The first, where the strands were placed consecutively in an ABCDEFG "up-down" pattern, found its closest analog in the retinol binding protein.²⁷ Including a single Greek key element in the fold approximated the fold found in pseudoazurin (ABEDCFG).²⁸ To make the analogy to pseudoazurin "work," the first strand of pseudoazurin was ignored, and a strand was moved from one sheet in the β -sandwich to the other (the "modified pseudoazurin fold").⁶ The third alternative fold (ABCDGFE) had a topology similar to that found in the pleckstrin homology.²⁹ The "modified pseudoazurin" fold (ABEDCFG) turned out to be the correct topology for the fold of first seven B-strands in synaptotagmin.⁹

The fact that the complexity of the tertiary structure modeling problem was reduced in this case from hundreds of thousands of possible topologies to just three, with the correct topology contained within the three, makes several points. First, the misassignment of the final strand in the fold as a helix did not interfere with the assignment of the tertiary fold. This was undoubtedly due to the fact that the misassigned secondary structural unit was not at the core of the fold (see above), and also because it came at the end of the domain, not at the middle. Of course, the fact that the segment was misassigned and that it is not a core segment is not independent, as noted above.

Second, we were unable to identify the correct tertiary fold from among these three alternatives. Efforts were made to identify hypotheses to try to create a preference for one of the three folding topologies. We noted that these hypotheses were little more than speculative efforts to use "a bona fide prediction opportunity to test some unorthodox ideas." Nevertheless, it is important to note that these unorthodox ideas failed. We preferred the pleckstrin homology domain fold over the pseudoazurin fold; as noted above, synaptotagmin in fact adopts (in its first 7 strands) a modified pseudoazurin topology. Application of a threading program might have been a better approach for identifying the correct fold from among the three preferred topologies.³⁰

Another noteworthy aspect of the synaptotagmin prediction is its model built for the calcium-binding active site. In the prediction, Asp-48 (Asp-178 in the synaptotagmin numbering), Asp-104 (Asp-230), Asp-106 (Asp-232), and Glu-81 (Glu-208) were predicted to form a calcium-binding site. In making this prediction, only patterns of conservation and variation within the multiple alignment were considered, together with the knowledge that synaptotagmin most likely contained a calcium-binding site. This model proved to be a good representation for the putative calcium binding active site in synaptotagmin, which is built from residues Asp-48, Asp-104, and Asp-106 (synaptotagmin numbering 178, 230, and 232). Thus, the active site assignments made in the prediction were useful guides for predicting tertiary structure, as they were in the phospho- β -galactosidase prediction.

WHAT HAVE WE LEARNED?

These two predictions add to the dozen or so remarkably accurate bona fide predictions that have now been made in many laboratories using aligned homologous protein sequences.¹⁵ As the collection of examples builds, much is learned about both the specific heuristics used to assign structural features of a protein, but also about protein folding in general. It seems that it is possible to obtain reasonably reliable predictions of surface and interior assignments from a set of aligned homologous protein sequences, and to use this tertiary structural information, together with a rule-based approach based for internal segments to obtain a plausible model of secondary structure, at least for core elements. These have been sufficient, at least in these cases, to narrow the number of possible folding topologies to a small number. In combination with threading tools. the number of tertiary models for the core can be reduced to one (for phospho- β -galactosidase). In other cases, additional tools, not yet developed, are needed to identify a unique prediction from a small set of alternative structures. We look forward to the next contest 2 years hence.

REFERENCES

 Hengstenberg, W., Kohlbrecher, D., Witt, E., Kruse, R., Christiansen, I., Peters, D., Vonstrandmann, R.P., Stadtler, P., Koch, B. Structure and function of proteins of the phosphotransferase system and of 6-phospho-beta-ga-

452

lactosidases in Gram-positive bacteria. FEMS Microbiol. Rev. 12:149-164, 1993.

- Gerloff, D.L., Benner, S.A. A consensus prediction of the secondary structure for the 6-phospho-β-D-galactosidase superfamily. Proteins 21:273-281, 1995.
 Benner, S.A., Jenny, T.F., Cohen, M.A., Gonnet, G.H. Pre-
- Benner, S.A., Jenny, T.F., Cohen, M.A., Gonnet, G.H. Predicting the conformation of proteins from sequences. Progress and future progress. Adv. Enzyme Reg. 34:269-353, 1994.
- Perin, M.S., Fried, V.A., Mignery, G.A., Jahn, R., Suedhof, T.C. Phospholipid binding by a synaptic vesicle protein homologous to the regulatory region of protein kinase C. Nature (London) 345:260-263, 1990.
- 5. Pevsner, A., Scheller, R. H. Mechanisms of vesicle docking and fusion. Insights from the nervous system. Curr. Opin. Cell Biol. 6:555-560, 1994.
- Gerloff, D.L., Chelvanayagam, G., Benner, S.A. A predicted consensus structure for the protein kinase C2 homology (C2H) domain, the repeating unit of synaptotagmin. Proteins 22:299-310, 1995.
- Cohen, F.E., Abarbanel, R.M., Kuntz, I.D., Fletterick, R. J. Secondary structure assignment for alpha-beta-proteins by a combinatorial approach. Biochemistry 22:4894-4904, 1983.
- Defay, T., Cohen, F.E. Evaluation of current techniques for ab initio protein structure prediction. Proteins 23:431-447, 1995.
- Sutton, R.B., Davletov, B.A., Berghuis, A.M., Südhof, T.C., Sprang, S.R. Structure of the first C2 domain of synaptotagmin I: A novel Ca²⁺ phospholipid binding fold. Cell 80:929-938, 1995.
- 80:929-938, 1995.
 Cohen, B.I., Presnell, S.R., Cohen, F.E. Origins of structural diversity within sequentially identical hexapeptides. Protein. Sci. 2:2134-2145, 1993.
 Benner, S.A. Patterns of divergence in homologous pro-
- Benner, S.A. Patterns of divergence in homologous proteins as indicators of tertiary and quaternary structure. Adv. Enzyme Reg. 31:121–181, 1989.
- 12. Benner, Š.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.
- 13. Benner, S.A. Predicting de novo the folded structure of proteins. Curr. Opin. Struct. Biol. 2:402-412, 1992.
- 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.
- 15. Benner, S.A., Gerloff, D.L., Jenny, T.F. Predicting protein crystal structures. Science 265:1642-1644, 1994.

- 16. Lim, V.I. Algorithms for prediction α -helices and β -structural regions in globular proteins. J. Mol. Biol. 88:873-894, 1974.
- Gerloff, D.L., Benner, S.A. Predicting the conformation of proteins: Man versus machine. FEBS Lett. 325:29-33, 1993.
- 18. Benner, S.A., Cohen, M.A., Gerloff, D.L. A predicted secondary structure for the Src homology domain 3. J. Mol. Biol. 229:295-305, 1993.
- Chothia, C., Lesk, A. The relation between the divergence of sequence and structure in proteins. EMBO J. 5:823-826, 1986.
- Zhu, Z.Y. A new approach to the evaluation of protein secondary structure predictions at the level of the elements of secondary structure. Prot. Eng. 8:103-108, 1995.
- Jenny, T.F., Benner, S.A. Evaluating predictions of secondary structure in proteins. Biochem. Biophys. Res. Commun. 200:149-155, 1994.
- Bryant, S.H., Lawrence, C.E. An empirical energy function for threading protein, sequence through the folding motif proteins. 16:92-112, 1993.
- Farber G.K., Petsko, G. The evolution of α/β barrel enzymes. Trends Biochem. Sci. 15:228-234, 1990.
- 24. Crawford, I.P., Niermann, T., Kirschner, K. Prediction of secondary structure by evolutionary comparison: Application to the α subunit of tryptophan synthase. Protein 2:118-129, 1987.
- Garnier, J., Robson, B. The GOR method for predicting secondary structures in proteins. In: "Prediction of Protein Structure and the Principles of Protein Conformation." G.D. Fasman, (ed.). Plenum, New York: 1989: 417-465.
- Woolfson, D.N., Evans, P.A., Hutchinson, E.G., Thornton, J.M. Topological and stereochemical restrictions in betasandwich protein structures. Prot. Eng. 6:461-470, 1993.
- Zanotti, G., Berni, R., Monaco, H.L. Crystal structure of liganded and unliganded forms of bovine plasma retinol binding protein. J. Biol. Chem. 268:10727-10738, 1993.
- Petratos, K., Banner, D.W., Beppu, T., Wilson, K.S., Tsernoglou, D. The crystal structure of pseudoazurin from Alkaligenes faecalis SS-6 determined at 2.9 Å resolution. FEBS Lett. 218:209-214, 1987.
- Yoon, H.S., Hajduk, P.J., Petros, A.M., Olejniczak, E.T., Meadows, R.P., Fesik, S.W. Solution structure of a pleckstrin-homology domain. Nature (London) 369:672-675, 1994.
- Jones, D.T., Taylor, W.R., Thornton, J.M.A new approach to protein fold recognition. Nature (London) 358:86-89, 1992.