

## Structure-based statistical analysis of transmembrane helices.

Carlos Baeza<sup>1</sup>, Marc A. Marti-Renom<sup>2,\*</sup> and Ismael Mingarro<sup>1,\*</sup>

1. Departament de Bioquímica i Biologia Molecular. Universitat de València.

Burjassot, Spain.

2. Structural Genomics Group. National Center for Genomic Analysis (CNAG).

Barcelona, Spain.

\* Corresponding authors:

Marc A. Marti-Renom

Ismael Mingarro

Phone: +34 934 033 743

Phone: +34 963 543 796

Fax: +34 934 037 279

Fax: +34 963 544 635

E-mail: [mmarti@cpb.uv.es](mailto:mmarti@cpb.uv.es)

E-mail: [Ismael.Mingarro@uv.es](mailto:Ismael.Mingarro@uv.es)

**Version:** 3/21/2012

## Abstract

Recent advances in high-resolution structure determination of membrane proteins enable now the analysis of the main features of amino acids in transmembrane (TM) segments in comparison with amino acids in water-soluble helices. In this work, we introduced a large-scale analysis of amino acid propensities using a data set of 170 structures of integral membrane proteins obtained from MPTopo database and 930 structures of water-soluble helical proteins obtained from the Protein Data Bank. Large hydrophobic residues (Leu, Val, Ile and Phe) plus Gly had a clear preference for TM helices, while polar residues (Glu, Lys, Asp, Arg and Gln) were less frequent in this type of helices. The distribution of residues along the TM helices was also examined. As expected, hydrophobic and slightly polar amino acids are commonly found in the hydrophobic core of the membrane, while aromatic (Trp and Tyr) and Pro together with hydrophilic (Asn, His, and Gln) residues are frequent in the interface regions. Charged residues also have statistically preferred locations avoiding the hydrophobic core of the membrane, but while acidic residues are frequently found at both the cytoplasmic and extra-cytoplasmic interfaces, basic residues cluster at the cytoplasmic interface. These results strongly support the experimentally demonstrated biased distribution of positively charged residues (that is, the so-called the positive-inside rule) with structural data.

## Keywords

Membrane protein; transmembrane helices; amino acid distribution; statistical analysis.

## Introduction

1  
2  
3 Although helical membrane proteins represent about one fourth of all proteins  
4 in living organisms (Wallin & Heijne, 1998), the rules governing its folding are  
5 still not completely established. The hydrophobic effect is a dominant driving  
6 force to the folding of water-soluble proteins, but its contribution to the folding  
7 of membrane proteins is further more complex given that these proteins “live”  
8 in a biophysical environment –the membrane–, which is clearly different from  
9 the aqueous media. The cell membrane is a very heterogeneous media,  
10 composed mainly of phospholipids that are self-organized in two leaflets  
11 giving rise to the formation of a bilayer. The hydrocarbon core is the  
12 hydrophobic part of the membrane, covering approximately 30 Å. The polar  
13 head groups of the phospholipids define the lipid/water interphase and add  
14 approximately 15 Å to the thickness of each leaflet (White & Wimley, 1999). It  
15 is in this complex environment in which membrane proteins have to fold into  
16 their native conformations.  
17

18  
19 The hydrocarbon core of the biological membranes and the interior of folded  
20 water-soluble proteins are hydrophobic. In such a hydrophobic environment,  
21 the polarity of the polypeptide backbone is energetically unfavorable. Thus, in  
22 protein structures, nearly all the polar groups of the peptide bond (carbonyl  
23 and amide groups) tend to hydrogen bond with one another, leading to  
24 secondary structure that stabilizes the folded state. Alpha-helices are the  
25 commonest secondary structural elements found in water-soluble as well as in  
26 membrane protein structures. However, the distribution of the helices in these  
27 two groups of proteins is very different. While helices in water-soluble  
28 proteins can be exposed to both the hydrophobic core and the water-  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 accessible surface, transmembrane (TM) helices in membrane proteins are  
2 surrounded by a hydrophobic lipid phase where water is essentially absent.  
3  
4 Therefore, for the structural stabilization of helical membrane proteins that  
5 reside in this apolar (low dielectric) environment, hydrogen bonding and van  
6 der Waals packing forces have an increased importance.  
7  
8  
9

10  
11 Although the great majority of membrane proteins integrate into biological  
12 membranes through the translocon (see for a recent review (Martínez-Gil et al,  
13 2011)), our current biophysical understanding of its folding and function is  
14 hampered by the scarcity of structural information. Fortunately, the number of  
15 high-resolution structures of membrane proteins has increased exponentially  
16 in the last years (White, 2004; 2009). Consequently, a new statistical survey  
17 of TM helices properties is timely.  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29

30 In this paper, we revisit the differences between helices from water-soluble  
31 proteins and TM helices in terms of length and amino acid composition. In  
32 addition, we analyze the distribution of amino acid residues in TM segments,  
33 which have to energetically accommodate into the highly heterogeneous  
34 media of biological membranes by interacting favorably with its local  
35 environment. The present study involved 170 helical membrane proteins with  
36 known three-dimensional structure and topology, containing a total of 792 TM  
37 segments and compared with 7,348 helices from 930 water-soluble protein  
38 structures. About half of all amino acids are randomly distributed when  
39 allocated into the membrane, but the rest show a strong correlation for  
40 residue positions along the TM regions.  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

## Methods

### *Helix data sets*

Two data sets of water-soluble and TM helices were obtained from the Protein Data Bank (PDB) (Berman et al, 2000) and the MPTOPO databases (Jayasinghe et al, 2001b), respectively.

First, a total of 4,405 structural chains deposited in the PDB (as of November 17<sup>th</sup>, 2011) that passed the following criteria were selected: (i) their total secondary structure had more than 60% of  $\alpha$ -helices and no  $\beta$ -strands; (ii) their crystallographic resolution was 2.0 Å or higher; and (iii) the word *MEMBRANE* did not appear in the “TITLE” nor the “DESCRIPTION” fields of the PDB file. Furthermore, to remove redundancy, the 4,405 chain sequences were compared to each other with the *cd-hit* program (Huang et al, 2010) and pairs resulting in sequence alignments with 80% or higher identity were discarded. The final set of 930 non-redundant PDB chains was parsed to identify a total of 7,348 helices from “HELIX” fields of each PDB chain entry. Thus, the data set of water-soluble helices contained 930 non-redundant and high-resolution protein structures, 7,348  $\alpha$ -helices and 108,277 amino acids.

Second, all  $\alpha$ -helical membrane proteins deposited in the MPTOPO database (last updated on January 19<sup>th</sup>, 2010) (Jayasinghe et al, 2001b), and thus with known membrane insertion topology, were selected. The initial set was further filtered by: (i) removing any entry of unknown structure as based on the MPTOPO entry classification (*i.e.*, keeping only entries described as “3D\_helix” and “1D\_helix”); (ii) removing redundant pairs at 80% sequence identity by applying the *cd-hit* program (Huang et al, 2010). The final data set

of TM helices contained 170 non-redundant structures, 837 TM helices, and 20,079 amino acids. Furthermore, to properly analyze the amino acid propensities in single membrane spanning TM helices, we discarded any helix shorter than 17 amino acids or larger than 38 amino acids. The resulting TM data subset contained 792 TM helices, and 19,356 amino acids.

### *Amino acid propensity measures*

We calculated three different amino acid measures: (i) probability and percent, (ii) Odds, and (iii) LogOdds. The probability ( $p_i$ ) of an amino acid  $i$  is defined as:

$$p_i = \frac{n_i}{N}$$

where  $i$  is the amino acid type (one of the 20 amino acids),  $n_i$  is the observation count of the amino acid  $i$ , and  $N$  is all amino acids in the data set.

Similarly, the percent of a given amino acid  $i$  is defined as its probability multiplied by 100. The Odds ( $O_i$ ) of an amino acid  $i$  is defined as:

$$O_i = \frac{p_{i,c}}{(1 - p_{i,c})} \bigg/ \frac{p_{i,r}}{(1 - p_{i,r})}$$

where  $p_{i,c}$  is the probability of the amino acid  $i$  in the class  $c$  (for example, TM helix) and  $p_{i,r}$  is the probability of the amino acid  $i$  in the class  $r$  (for example, water-soluble helix). Similarly, the LogOdds of a given amino acid  $i$  is defined as the logarithm in base 10 of its Odds. Briefly, Odds higher than 1 (or positive LogOdds) indicate over-occurrence of the amino acid type in the class. Odds smaller than 1 (or negative LogOdds) indicate under-representation of the amino acid type in the class.

## Results and Discussion

### *Helix length in membrane and water-soluble proteins*

Length distributions for helices found in high-resolution structures deposited in PDB (Berman et al, 2000) are very different for TM and water-soluble proteins (Fig. 1).

Helices in TM proteins are in average 24.0 ( $\pm$  5.6) amino acid residues long, this result slightly differs from previous data obtained using databases with 45 (Bowie, 1997) and 129 (Ulmschneider & Sansom, 2001) TM helices, where average helix length was 26.4 and 27.1 amino acid residues, respectively. As the translation per residue in a canonical helix is 1.5Å, a stretch of about 20 consecutive hydrophobic residues can span the 30 Å of the hydrocarbon core of biological membranes. Indeed, the more prevalent (~12%) length for TM helices in our data set was 21 residues (Fig. 1). Longer helices can span the bilayer with a concomitant tilting of the helix axis respect to the membrane plane. Other options are also feasible ranging from lipid accommodation till polypeptide backbone deformation (Holt & Killian, 2009).

Helices from water-soluble proteins have an average length of 14.7 ( $\pm$ 8.7) residues, which agrees with previous studies where the more prevalent helix length was 10-11 residues long (Engel & DeGrado, 2004; Pal et al, 2003). The reduced length for helices from water-soluble proteins is due to the absence of the restrictions imposed by the low dielectric constant at the hydrocarbon core of biological membranes, which forces the polypeptide backbone to adopt on average larger secondary structures.

### *Amino acid composition of $\alpha$ -helices*

1  
2  
3 The amino acid composition for both, TM and water-soluble helices, have  
4  
5 been examined (Fig. 2). TM helices of lengths between 17 and 38 residues  
6  
7 were selected from the MPTOPO database (Jayasinghe et al, 2001b), which  
8  
9 included helical segments that do completely span the hydrophobic core of  
10  
11 the membrane. TM helices shorter than 17 residues as well as larger than 38  
12  
13 residues were excluded since they may not cross entirely the membrane (Fig.  
14  
15 1 inset a) or may contain segments parallel to the membrane (Fig. 1 inset b).  
16  
17 Note that in the case of water-soluble helices all lengths were included in our  
18  
19 analysis because no restrictions in terms of length can be assumed for water-  
20  
21 soluble proteins in an aqueous milieu.  
22  
23  
24  
25  
26

27  
28 As expected, hydrophobic residues Leu, Ala, Val and Ile constitute the bulk of  
29  
30 the amino acids in the TM region accounting for almost half (47.0%) of all  
31  
32 residues. Similarly, these residues are also frequently found in helices of  
33  
34 water-soluble proteins (34.1%). However, there are, as noted previously  
35  
36 using smaller datasets (Bywater et al, 2001), differences in composition of the  
37  
38 two types of helices. Despite sharing the same structural features, the  
39  
40 differences between the two types of helices are reflected by their preferential  
41  
42 occurrences measured by the logarithm of the Odds of finding a given amino  
43  
44 acid in a TM helix with respect to its frequency in a water-soluble helix (Fig. 2  
45  
46 bottom panel). For example, while charged and polar residues are much  
47  
48 more frequently found in helices from water-soluble proteins, Trp, Gly and  
49  
50 Phe have higher propensities in TM helices. Interestingly, in contrast to their  
51  
52 conformational preferences in water, the helical propensities of residues such  
53  
54 as Val, Ile, Phe and Met are notably increased in the membrane environment,  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65



1 where it has been suggested that their helical proclivity is primarily governed  
2 by their side chain hydrophobicity and by the hydrophobicity of the local  
3 polypeptide region in which the residues reside span the membrane (Li &  
4 Deber, 1994). Significantly, Gly and Pro are more frequent in TM helices  
5 relative to water-soluble helices. Although commonly considered as 'helix  
6 breakers' it has been reported that Gly residues occur frequently in TM helix-  
7 helix interactions, especially in association with  $\beta$ -branched residues at  
8 neighboring positions (Senes et al, 2000), and that Pro, in addition to its role  
9 in signal transduction and gating across the membrane, may also play a  
10 significant role in these processes (Orzáez et al, 2004).

11 A comparison of the amino acid frequency between TM and water-soluble  
12 helices confirmed that strongly polar residues (Glu, Lys, Asp, Arg, and Gln)  
13 are more prevalent in water-soluble helices (Fig. 3). These residues  
14 constitute only 8.2 % of the residues within TM helices compared to 30.9 % in  
15 water-soluble helices. Despite their lower presence, polar residues are  
16 evolutionary conserved in TM proteins, which has been partially explained by  
17 their tendency to be buried in the protein interior and also in many cases due  
18 to their direct involvement in the function of the protein (Illergård et al, 2011).  
19 Conversely, hydrophobic amino acids (Leu, Val, Ile, Gly, and Phe) are over-  
20 represented in TM helices (Fig. 3). Interestingly, Ala although being the  
21 second more abundant residue in TM helices (Fig. 2), it is not over-  
22 represented in this type of helices likely because its higher helical propensity  
23 in aqueous (Blaber et al, 1993) compared to membrane-mimetic  
24 environments (Li & Deber, 1994). In fact, both biological (Nilsson et al, 2003;  
25 Hessa et al, 2005) and biophysical (Jayasinghe et al, 2001a) measurements  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 have poised Ala at the threshold between those amino acids that promote  
2 membrane integration of TM helices and those residues that preclude  
3 membrane insertion.  
4  
5

6  
7  
8 *Position dependent distribution of amino acid residues in TM helices*  
9

10 A comparison of the amino acid frequency at different positions in a TM  
11 segment, taking as reference the TM center, confirmed that about half of the  
12 natural amino acid residues have similar distributions at positive positions  
13 (towards inside the cell) than at negative positions (towards outside the cell)  
14 (Fig. 4). It was found that not only the strongly hydrophobic residues but also  
15 Gly and the hydroxylated residues Ser and Thr are equally distributed along  
16 the hydrophobic core of the membrane. It is important to note that Gly is a  
17 residue type that is normally regarded as being conducive to turn (Williams et  
18 al, 1987), yet it is a common residue in TM helices (Fig. 2). There are  
19 important folding reasons for incorporating Gly into TM helices. The absence  
20 of side-chain of the Gly allows for bulkier groups to be accommodated close  
21 to the polypeptide backbone of the TM helices. This might be important for  
22 intramolecular helix-helix packing, for homo-oligomerization, or for recognition  
23 of other membrane proteins, among other factors. Indeed, it has been  
24 observed that Gly has the highest overall packing value in membrane proteins  
25 (Eilers et al, 2002). Ser or Thr residues within TM helices participate in  
26 hydrogen-bonding networks through hydrogen bond linking of the side chain  
27 oxygen atom to acceptor side chain or peptide bond groups. These effects,  
28 intimate packing (Gly) and hydrogen bonding (Ser and Thr), can be relevant  
29 at any position along the TM region, which would explain the absence of  
30 position preference for these residues in TM helices. Met or Cys are also  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 frequent at different locations within the hydrophobic core, but a relative  
2 prevalence can be observed in a region that would correspond with the initial  
3 portion of the polar headgroups of the phospholipids, consistent with the  
4 slightly amphipathic nature of these residues and in agreement with its  
5 distribution in the lipid bilayer recently obtained from molecular dynamics  
6 simulation (MacCallum et al, 2008).  
7  
8  
9  
10  
11  
12  
13

14 While Phe has a flat distribution in TM helices, behaving as a hydrophobic  
15 residue, Trp, Tyr and Pro residues are distributed in a biased manner: they  
16 are found preferentially at the ends of the bilayer (*i.e.* at the interface between  
17 the hydrophobic core of the bilayer and the bulk water). At this location,  
18 aromatic residues may serve as anchors for the TM helices into the  
19 membrane. In fact, Trp and Tyr positioned 7 to 9 residues away from the  
20 center of a TM segments result in a reduction in free energy (Hessa et al,  
21 2007), which nicely correlates with the present statistical distribution from  
22 three-dimensional structures (Fig. 4). The biophysical reason for the  
23 observed distribution of Trp and Tyr residues could rely on the relatively  
24 amphipathic nature of their side chains, which can form hydrogen bonds as  
25 well as exhibit hydrophobic character. Actually, this preferred location has  
26 previously been observed not only for  $\alpha$ -helical but also  $\beta$ -barrel membrane  
27 proteins (Ulmschneider & Sansom, 2001). A similar distribution is observed  
28 for Pro residues, although an increased presence is detectable towards the  
29 center of the bilayer, which can be associated with the fundamental and  
30 subtle role that Pro residues play in the dynamics, structure and function of  
31 many membrane proteins by inducing the formation of molecular hinges  
32 (Cordes et al, 2002). Indeed, thirteen TM helices with known structure have a  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 Pro residue at the 0 position, which in all cases results in a kink in the helix.  
2 Nevertheless, it should be noted that the interfacial preference of these three  
3 residues is somehow more pronounced at the non-cytoplasmic interface. This  
4 was also observed in the case of aromatic residues (Trp and Tyr) in a  
5 membrane protein prediction analysis using sequence information from 107  
6 genomes (Nilsson et al, 2005).  
7  
8  
9  
10  
11  
12  
13

14 The distribution pattern for Asn, His and Gln, corresponds to an interfacial  
15 preference close to the end of the TM regions, which is consistent with the  
16 amphipathic nature of these molecules. This pattern was previously reported  
17 for His residues (Ulmschneider & Sansom, 2001), which is in good agreement  
18 with our data. Interestingly, in more recent studies using computer  
19 simulations, it has been noted that small molecule analogs of Asn (MacCallum  
20 et al, 2008) and Asn, His, and Gln (Johansson & Lindahl, 2007) result in an  
21 energetic minimum for partition into model lipid bilayers.  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33

34 Since the energetic cost of inserting an ionizable group in the hydrophobic  
35 environment of the membrane is very high (White & Wimley, 1999), charged  
36 amino acids should generally be excluded from the hydrophobic core of the  
37 TM helices. Interestingly, nearly all membrane proteins with six or more  
38 predicted TM helices contain at least one ionizable residue (Arkin & Brunger,  
39 1998). However, charged amino acids consistently clustered at the TM  
40 flanking regions (Fig. 4). For example, acidic (Asp and Glu) residues result in  
41 an increased distribution at both cytoplasmic and extra-cytoplasmic side of the  
42 membrane, although with some prevalence for the cytoplasmic region.  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 with the positive-inside rule (Heijne, 1992). Moreover, it has been  
2 experimentally demonstrated that basic residues act as stronger topological  
3 signals compared to acidic residues (Nilsson & Heijne, 1990; Saurí et al,  
4 2009), which is reflected by their different statistical preferences on either end  
5 of the TM segments. Nevertheless, when considered globally, charged  
6 residues cluster preferentially near the cytoplasmic end of the TM segments  
7 (Fig. 5, orange line). This effect was already noted in a previous structure-  
8 based analysis that included a lower number of structures available at the  
9 time (Ulmschneider et al, 2005). On the contrary, although polar residues  
10 (Gln, His, and Asn) mimic the distribution pattern of charged residues avoiding  
11 the more hydrophobic region of the bilayer, they show a preference for the  
12 extra-cytoplasmic region (Fig. 5). Trp, Tyr and Pro are more abundant about 8  
13 to 9 residue positions away from the center of the membrane, that is, within  
14 the interface region, but with some bias toward the extra-cytoplasmic interface.  
15 The rest of natural amino acids are more abundant at the center of the bilayer,  
16 within 7 amino acid positions on both sides of the membrane normal, but they  
17 are also very frequently found beyond this boundary as noted by their overall  
18 proximity to the Odd value of 1 for positions >10 on both sides of the center of  
19 the membrane (Fig. 5). Interestingly, the amino acid distribution patterns at  
20 both interface regions are slightly different. There is a sharper transition from  
21 mainly hydrophobic to charged, polar and aromatic residues at the  
22 cytoplasmic side of the membrane (positions 6 to 8) compared to that at the  
23 extra-cytoplasmic side (positions -5 to -9). The different lipid composition  
24 between the two lipid leaflets in biological membranes and the strong  
25 electrochemical potential over the prokaryotic inner cell membranes can exert

1 an important effect, which may be reflected by this difference. For instance, it  
2 has been recently reported an asymmetry in the distribution of amino acid  
3 residues within TM segments from plasma membrane proteins (Sharpe et al,  
4 2010), which has been attributed to an asymmetry in the state of lipid order in  
5 the membrane. Such an asymmetry is likely due to the enrichment of lipids  
6 such as sterols and sphingolipids in the extra-cytoplasmic leaflet, where a  
7 more gradual amino acid distribution can be expected.  
8  
9

10 Finally, we analyzed and plotted the odd ratio for each amino acid in three  
11 regions in a membrane, that is, taking the hydrophobic TM region as the  
12 central 19 positions ( $\sim 30\text{\AA}$ ) and 9 residue positions ( $\sim 15\text{\AA}$ ) on both sides as  
13 the extra-cytoplasmic (from -10 to -18 residues) and cytoplasmic (from 10 to  
14 18) flanking regions (Fig. 6). Hydrophobic amino acids (blue colored)  
15 populated preferentially the hydrophobic center. However, this trend is not  
16 observed for the more prevalent residues in TM segments (for example Leu,  
17 Fig. 2), which are also frequently found at the flanking regions. Trp, Tyr, and  
18 Pro (green) have a minor increase for the extra-cytoplasmic flanking region.  
19 The absence of higher differences for the distribution of these residues is  
20 probably due to their precise location at the interface between the  
21 hydrophobic core and the flanking hydrophilic environment. Polar (orange)  
22 residues (Gln, His, and Asn) have a preference for both flanking regions since  
23 they are energetically unfavorable within the membrane core. These residues  
24 do not ionize at the physiological pH and are able to donate and accept  
25 hydrogen bonds simultaneously. Such an effect translates into a higher  
26 preference of Gln, His and Asn for the rich hydrogen bond network  
27 environment of the interface. Charged residues (red) are underrepresented at  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 the hydrophobic core and resulted in preferences for the cytoplasmic flanking  
2 region being acidic residues more prevalent at the extra-cytoplasmic flanking  
3 region. Furthermore, basic residues are strong topological determinants that  
4 heavily populate the cytoplasmic flanking region. The effect of positively  
5 charged residues located near the cytoplasmic end of hydrophobic segments  
6 has been in fact estimated to be approximately -0.5 kcal/mol to the apparent  
7 free energy of membrane insertion (Lerch-Bader et al, 2008). This energetic  
8 contribution can be extremely relevant to precisely anchor hydrophobic  
9 regions into biological membranes.  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24

### 25 **Concluding remarks**

26  
27 We have compared the length and the amino acid composition of helices in  
28 TM and water-soluble proteins. Overall, significant differences are present in  
29 both types of proteins, which may be attributed to the biophysical differences  
30 between the two environments in which they fold. First, TM helices adapt  
31 their length to the dimensions and constraints of biological membranes, while  
32 water-soluble helices are statistically shorter since they do not have to satisfy  
33 the demanding restrictions imposed by the complexity of the membrane  
34 environment. Second, the observed differences highlight that in the lipid  
35 bilayer, which environment forces secondary structure formation, amino acid  
36 side chain hydrophobicity prevails to helicity. Accordingly, aliphatic residues  
37 with a reduced helical propensity (Val, Ile, Gly, and Phe) are abundant in TM  
38 helices, while polar residues (Glu, Lys, and Arg) with high helical propensity  
39 are consistently less frequent in TM helices. Third, half of the natural amino  
40 acid residues are equally distributed along the TM helices, whilst aromatic,  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 polar and charged residues plus Pro are biased toward the ends of the TM  
2 helices. Fourth, as previously observed, the distribution of charged residues  
3 was asymmetric occurring more frequently on the cytoplasmic side of the  
4 membrane, which causes a net charge unevenness on both sides of the  
5 membrane. In addition to this asymmetry, Trp, Tyr and Pro residues were  
6 found to be more frequent at the extra-cytoplasmic interface of the membrane  
7 and the polar residues (Gln, His, and Asn) at the extra-cytoplasmic flanking  
8 region of the TM helices. Fifth, transitions between the different types of  
9 residues at the ends of the hydrophobic core occur in a more defined region  
10 at the cytoplasmic side than at the extra-cytoplasmic face, likely reflecting the  
11 differences in lipids composition on both leaflets of biological membranes.  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26

27 The conclusions on TM helix architecture described here should prove useful  
28 for constructing models of membrane proteins with desired properties, which  
29 could help filling in some of the many gaps in the field.  
30  
31  
32  
33  
34  
35  
36  
37

### 38 **Acknowledgments**

39  
40  
41 This work was supported by grants BFU2009-08401 (to I.M.) and BFU2010-  
42 19310 (to M.A.M-R.) from the Spanish Ministry of Science and Innovation  
43 (MICINN, ERDF supported by the European Union), as well as by  
44 PROMETEO/2010/005 and ACOMP/2012/226 (to I.M.) and ACOMP/2011/048  
45 (to M.A.M-R.) from the Generalitat Valenciana. C.B. was recipient of a  
46 predoctoral FPI fellowship from the MICINN.  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65



## References

- 1  
2 Arkin IT & Brunger AT (1998) Statistical analysis of predicted transmembrane  
3 alpha-helices. *Biochim Biophys Acta* **1429**: 113–128  
4
- 5  
6 Berman HM, Berman HM, Westbrook J, Westbrook J, Feng Z, Feng Z,  
7 Gilliland G, Gilliland G, Bhat TN, Bhat TN, Weissig H, Weissig H,  
8 Shindyalov IN, Shindyalov IN, Bourne PE & Bourne PE (2000) The  
9 Protein Data Bank. *Nucleic Acids Res.* **28**: 235–242  
10
- 11  
12 Blaber M, Zhang XJ & Matthews BW (1993) Structural basis of amino acid  
13 alpha helix propensity. *Science* **260**: 1637–1640  
14
- 15  
16 Bowie JU (1997) Helix packing in membrane proteins. *J Mol Biol* **272**: 780–  
17 789  
18
- 19  
20 Bywater RP, Thomas D & Vriend G (2001) A sequence and structural study of  
21 transmembrane helices. *J. Comput. Aided Mol. Des.* **15**: 533–552  
22
- 23  
24 Cordes FS, Bright JN & Sansom MSP (2002) Proline-induced distortions of  
25 transmembrane helices. *J Mol Biol* **323**: 951–960  
26
- 27  
28 Eilers M, Patel AB, Liu W & Smith SO (2002) Comparison of helix interactions  
29 in membrane and soluble alpha-bundle proteins. *Biophys J* **82**: 2720–  
30 2736  
31
- 32  
33 Engel DE & DeGrado WF (2004) Amino acid propensities are position-  
34 dependent throughout the length of alpha-helices. *J Mol Biol* **337**: 1195–  
35 1205  
36
- 37  
38 Heijne von G (1992) Membrane protein structure prediction. Hydrophobicity  
39 analysis and the positive-inside rule. *J Mol Biol* **225**: 487–494  
40
- 41  
42 Hessa T, Kim H, Bihlmaier K, Lundin C, Boekel J, Andersson H, Nilsson I,  
43 White SH & Heijne von G (2005) Recognition of transmembrane helices  
44 by the endoplasmic reticulum translocon. *Nature* **433**: 377–381  
45
- 46  
47 Hessa T, Meindl-Beinker NM, Bernsel A, Kim H, Sato Y, Lerch-Bader M,  
48 Nilsson I, White SH & Heijne von G (2007) Molecular code for  
49 transmembrane-helix recognition by the Sec61 translocon. *Nature* **450**:  
50 1026–1030  
51
- 52  
53 Holt A & Killian JA (2009) Orientation and dynamics of transmembrane  
54 peptides: the power of simple models. *Eur Biophys J* **39**: 609–621  
55
- 56  
57 Huang Y, Huang Y, Niu B, Niu B, Gao Y, Gao Y, Fu L, Fu L, Li W & Li W  
58 (2010) CD-HIT Suite: a web server for clustering and comparing biological  
59 sequences. *Bioinformatics* **26**: 680–682 Available at:  
60 <http://eutils.ncbi.nlm.nih.gov/entrez/eutils/elink.fcgi?dbfrom=pubmed&id=20053844&retmode=ref&cmd=prlinks>  
61  
62  
63  
64  
65

membrane core evolutionary conserved? *Proteins* **79**: 79–91

Jayasinghe S, Hristova K & White SH (2001a) Energetics, stability, and prediction of transmembrane helices. *J Mol Biol* **312**: 927–934

Jayasinghe S, Jayasinghe S, Hristova K, Hristova K, White SH & White SH (2001b) MPtopo: A database of membrane protein topology. *Protein Sci* **10**: 455–458

Johansson ACV & Lindahl E (2007) Position-resolved free energy of solvation for amino acids in lipid membranes from molecular dynamics simulations. *Proteins* **70**: 1332–1344

Lerch-Bader M, Lundin C, Kim H, Nilsson I & Heijne von G (2008) Contribution of positively charged flanking residues to the insertion of transmembrane helices into the endoplasmic reticulum. *Proc Natl Acad Sci USA* **105**: 4127–4132

Li SC & Deber CM (1994) A measure of helical propensity for amino acids in membrane environments. *Nat Struct Biol* **1**: 558

Lomize MA, Pogozheva ID, Joo H, Mosberg HI, Lomize AL (2012) OPM database and PPM web server: resources for positioning of proteins in membranes. *Nucleic Acids Res.* **40**: 370-376

MacCallum JL, Bennett WFD & Tieleman DP (2008) Distribution of amino acids in a lipid bilayer from computer simulations. *Biophys J* **94**: 3393–3404

Martínez-Gil L, Saurí A, Marti-Renom MA & Mingarro I (2011) Membrane protein integration into the endoplasmic reticulum. *FEBS J* **278**: 3846–3858

Nilsson I & Heijne von G (1990) Fine-tuning the topology of a polytopic membrane protein: role of positively and negatively charged amino acids. *Cell* **62**: 1135–1141

Nilsson I, Johnson AE & Heijne von G (2003) How hydrophobic is alanine? *J Biol Chem* **278**: 29389–29393

Nilsson J, Persson B & Heijne von G (2005) Comparative analysis of amino acid distributions in integral membrane proteins from 107 genomes. *Proteins* **60**: 606–616

Orzáez M, Salgado J, Giménez-Giner A, Pérez-Payá E & Mingarro I (2004) Influence of proline residues in transmembrane helix packing. *J Mol Biol* **335**: 631–640

Pal L, Chakrabarti P & Basu G (2003) Sequence and structure patterns in proteins from an analysis of the shortest helices: implications for helix nucleation. *J Mol Biol* **326**: 273–291

- 1 Saurí A, Tamborero S, Martínez-Gil L, Johnson AE & Mingarro I (2009) Viral  
2 membrane protein topology is dictated by multiple determinants in its  
3 sequence. *J Mol Biol* **387**: 113–128
- 4 Senes A, Gerstein M & Engelman DM (2000) Statistical analysis of amino  
5 acid patterns in transmembrane helices: the GxxxG motif occurs  
6 frequently and in association with beta-branched residues at neighboring  
7 positions. *J Mol Biol* **296**: 921–936
- 8  
9
- 10 Sharpe HJ, Stevens TJ, Munro S (2010) A comprehensive comparison of  
11 transmembrane domains reveals organelle-specific properties. *Cell* **142**,  
12 158-169.
- 13  
14
- 15 Ulmschneider MB & Sansom MS (2001) Amino acid distributions in integral  
16 membrane protein structures. *Biochim Biophys Acta* **1512**: 1–14
- 17  
18
- 19 Ulmschneider MB, Sansom MSP & Di Nola A (2005) Properties of integral  
20 membrane protein structures: derivation of an implicit membrane potential.  
21 *Proteins* **59**: 252–265
- 22  
23
- 24 Wallin E & Heijne von G (1998) Genome-wide analysis of integral membrane  
25 proteins from eubacterial, archaean, and eukaryotic organisms. *Protein*  
26 *Sci* **7**: 1029–1038
- 27  
28
- 29 White SH (2004) The progress of membrane protein structure determination.  
30 *Protein Sci* **13**: 1948–1949
- 31  
32
- 33 White SH (2009) Biophysical dissection of membrane proteins. *Nature* **459**:  
34 344–346
- 35  
36
- 37 White SH & Wimley WC (1999) Membrane protein folding and stability:  
38 physical principles. *Annu Rev Biophys Biomol Struct* **28**: 319–365
- 39  
40
- 41 Williams RW, Chang A, Juretić D & Loughran S (1987) Secondary structure  
42 predictions and medium range interactions. *Biochim Biophys Acta* **916**:  
43 200–204
- 44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

## Figure Legends

1  
2  
3  
4 **Figure 1.** Length distributions for 837 TM and 7,348 water-soluble helices  
5 from a set of non-redundant proteins of known structure (see Methods).  
6 Transmembrane helices in blue (pale blue correspond to discarded lengths)  
7 and water-soluble helices in orange. (a) Example of a short 9 amino acid  
8 length helix in the CIC chloride channel from *E. coli* (1KPK entry in PDB).  
9 Membrane boundaries were obtained from the PPM Server (Lomize et al,  
10 2012). The selected membrane is shown in rainbow coloring from N- (blue) to  
11 C-terminal (red) ends. (b) Example of a large 43 amino acid length helix in the  
12 chicken cytochrome BC1 complex (1BCC entry in PDB), which N-terminus of  
13 the helix (blue) lies at the membrane/water interface. Representation as in  
14 inset (a).  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33

34 **Figure 2.** Amino acid type distribution from 792 TM and 7,348 water-soluble  
35 helices from a set of non-redundant proteins of known structure (see  
36 Methods). (Upper plot) Amino acid type distribution for TM helices in blue and  
37 for water-soluble helices in orange. (Lower plot) LogOdds values for  
38 comparing the relative abundance of each amino acid type in TM and water-  
39 soluble helices. Amino acid types are ordered by its LogOdds  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51

52 **Figure 3.** Amino acid type percentage comparison between TM and water-  
53 soluble helices. Blue colored amino acids are over represented (difference >  
54 3 % points) in TM helices compared to water-soluble helices. Orange colored  
55 amino acids are over represented (difference > 3 % points) in water-soluble  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 helices compared to TM helices. Dashed grey lines indicate a cut-off of 3 %  
2 difference points.  
3

4  
5  
6  
7 **Figure 4.** Amino acid type and position distribution in TM helices. Each  
8 amino acid type and their positioning in the TM helix is represented by their  
9 positional normalized Odds (that is, for each column the Odds are normalized  
10 to an average of zero and standard deviation of one). The amino acids are  
11 clustered based on their positional normalized Odds within the helices.  
12 Positively labeled positions refer to the cytoplasmic side of the membrane and  
13 its flanking region whilst negatively labeled positions refer to extra-cytoplasmic  
14 regions.  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30

31 **Figure 5.** Amino acid groups positional preferences in a membrane. Thin  
32 lines represent the positional Odds for each amino acid individually, whilst  
33 thick lines represent the averaged positional Odds for each group of amino  
34 acids obtained from Figure 4. Amino acid types are grouped as in the  
35 dendogram in Fig. 4. That is, charged residues (red, KRED), polar residues  
36 (orange, QHN), aromatic residues plus Pro (green, PYW), and the rest of  
37 residues (blue, CMTSGVFAIL).  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50

51 **Figure 6.** Amino acid location preferences in a membrane. Letter size is  
52 proportional to the odds (relative preference) of finding a given amino acid in  
53 the three regions in a membrane (*i.e.*, from top to bottom outer, membrane,  
54 and inner regions). Amino acids colored as in figure 5.  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65













