

Rapport : Assignment et détection **des parties transmembranaires d'une protéine**

INTRODUCTION :

Les protéines membranaires représentent environ un quart de la totalité des séquences protéiques issues des gènes humains et assurent des fonctions clés pour l'organisme, notamment dans la communication et le transport entre les cellules et le milieu extra-cellulaire. De ce fait, ces protéines représentent des cibles d'intérêts majeurs dans le développement de nouveaux médicaments qui cherchent à agir directement sur leur fonctionnement. Cependant, les structures des protéines membranaires sont encore mal résolues aujourd'hui, principalement car réaliser leur cristallisation en environnement aqueux représente un challenge technique important.

Les nouvelles technologies de biologie structurale ont néanmoins permis d'apporter de nouvelles informations sur les coordonnées atomiques de nombreuses protéines membranaires. Toutefois, il reste très difficile de déterminer la position de la bicouche lipidique par rapport à la structure de ces protéines en conditions naturelles, car les possibilités d'obtenir des données expérimentales sur les membranes sont limitées. Pour reconstituer l'orientation des protéines au niveau des bicouches lipidiques, il faut donc se reporter à des méthodes prédictives proposées par des algorithmes dédiés tels que TMDET, ANVIL ou encore MEMEBED. L'utilisation de tels méthodes permet d'aboutir à une annotation des structures des protéines membranaires et de leurs segments transmembranaires, ce qui contribue à alimenter et mettre à jour les bases de données qui leurs sont consacrées, à l'image de PDB_TM ou OPM.

Dans ce contexte, l'objectif du travail présenté ici était d'établir un outil qui permettrait de déterminer les zones transmembranaires d'une protéine. Celui-ci est disponible à l'adresse : https://github.com/flotep/Projet_assignment_parties_transmb

La méthode proposée repose sur une approche géométrique pour suggérer la position, l'inclinaison, et l'épaisseur optimale de la membrane par rapport aux segments d'une protéine sur la seule base des coordonnées atomiques et de l'accessibilité au solvant de ses Calphas.

MATERIEL ET METHODES

Calcul de la surface accessible au solvant avec le programme DSSP

L'outil proposé traite les fichiers d'entrées en format PDB de la façon suivante. D'abord, il lit le fichier de la protéine et le passe au programme DSSP (*Define Secondary Structure of Proteins*)

pour calculer l'accessibilité au solvant de chacun de ces résidus à partir des structures secondaires qui sont assignées en fonction des coordonnées atomiques. Parmi ces résidus ne sont considérés que ceux pouvant interagir avec la bicouche lipidique, dont la valeur relative de la surface accessible (Relative ASA) est supérieure à un seuil défini à hauteur de 20% minimum.

Extractions des coordonnées des Calpha et calcul du centre de masse.

Les résidus de la protéine conservés sont ensuite divisés de manière binaire en deux catégories : hydrophobes pour les acides aminés {Phe, Met, Ile, Leu, Val, Trp, Ala, Cys, Gly, Ser, His}, et hydrophiles pour les acides aminés {Arg, Asp, Lys, Glu, Asn, Gln, Pro, Thr, Tyr}. Pour chaque résidu, seules les coordonnées atomiques des Calphas sont récupérées à partir du fichier PDB d'origine. Le centre de masse de la protéine est alors calculé à partir de ces Calphas et leurs coordonnées sont transformées par rapport à ce centre qui correspond donc au point (0,0,0) de notre repère.

Détermination de vecteur normaux pour orienter la membrane.

A partir de ce centre de masse est ensuite générée une surface sphérique sur laquelle sont échantillonnés des dizaines de points de manière aléatoire. Chaque point permet de créer un vecteur normal qui définit une position initiale et une direction que va prendre le plan de la membrane, orthogonal au vecteur selon l'équation : $ax + by + cz + d = 0$. L'objectif étant de quadriller avec suffisamment de précision toutes les directions possibles.

Par la suite, la position du plan de la membrane va être déplacée par tranche de 2 Angstrom le long de chaque vecteur normal tant que le nombre de résidus hydrophobes présent dans le plan est supérieur à 0. Dans le cas contraire, on estime que le plan que l'on regarde est situé en dehors de la protéine.

Scoring pour trouver les meilleures positions de la membrane.

Le but de cette approche est de calculer grâce à la fonction objective, pour chaque position de la membrane, le rapport entre le nombre d'atomes Calphas issus de résidus hydrophobes (M) à l'intérieur (Mi) et à l'extérieur (Me) du plan, par rapport au nombre d'atomes Calphas issus de résidus hydrophiles (S) à l'intérieur (Si) et à l'extérieur (Se) de ce même plan, selon la formule présentée dans *Postic et al*, 2015 :

$$C = \frac{M_i \times S_e - S_i \times M_e}{\sqrt{(M_i + S_i) \times (M_i + M_e) \times (S_i + S_e) \times (S_e + M_e)}}$$

Les acides aminés M sont prédits comme étant interne à la membrane tandis que les acides aminés S sont attendus à l'extérieur de cette membrane. La formule utilisée comme méthode de scoring dans cet algorithme cherche donc à trouver le plan de bicouche lipidique qui maximise la valeur de C.

Epaississement des meilleures positions de la membrane pour trouver le plan optimal.

Une fois que les scores de tous les plans possibles de la membrane ont été calculées, on conserve la position qui maximise C dans chaque direction. Pour chacune de ces positions, l'épaisseur de la membrane (initialement égale à 14 Angstrom) est augmentée par tranche de 1 Angstrom tant que la valeur de C augmente.

La position avec l'épaisseur de membrane pour laquelle le score C est le plus élevé donnera donc le plan optimal de la bicouche lipidique pour la protéine étudiée, l'équation du plan de chaque couche pouvant être récupérée en donnée de sortie.

Visualisation de la membrane avec la structure de la protéine.

Les sorties de cet algorithme permettent également une visualisation directe du plan optimal de la membrane par rapport à la structure de la protéine. D'abord, les coordonnées de tous les atomes du fichier PDB d'origine sont transformées par rapport au centre de masse précédemment calculé. Ensuite, ce fichier est édité en ajoutant des hétéroatomes dont les coordonnées ont été générées sur chaque plan de la bicouche lipidique de manière aléatoire. Ainsi, deux « tapis » de molécules sont ajoutés à la structure de la protéine et permettent d'identifier la position de la membrane en ouvrant le nouveau fichier PDB avec un logiciel de visualisation tel que PyMOL.

RESULTATS ET DISCUSSION :

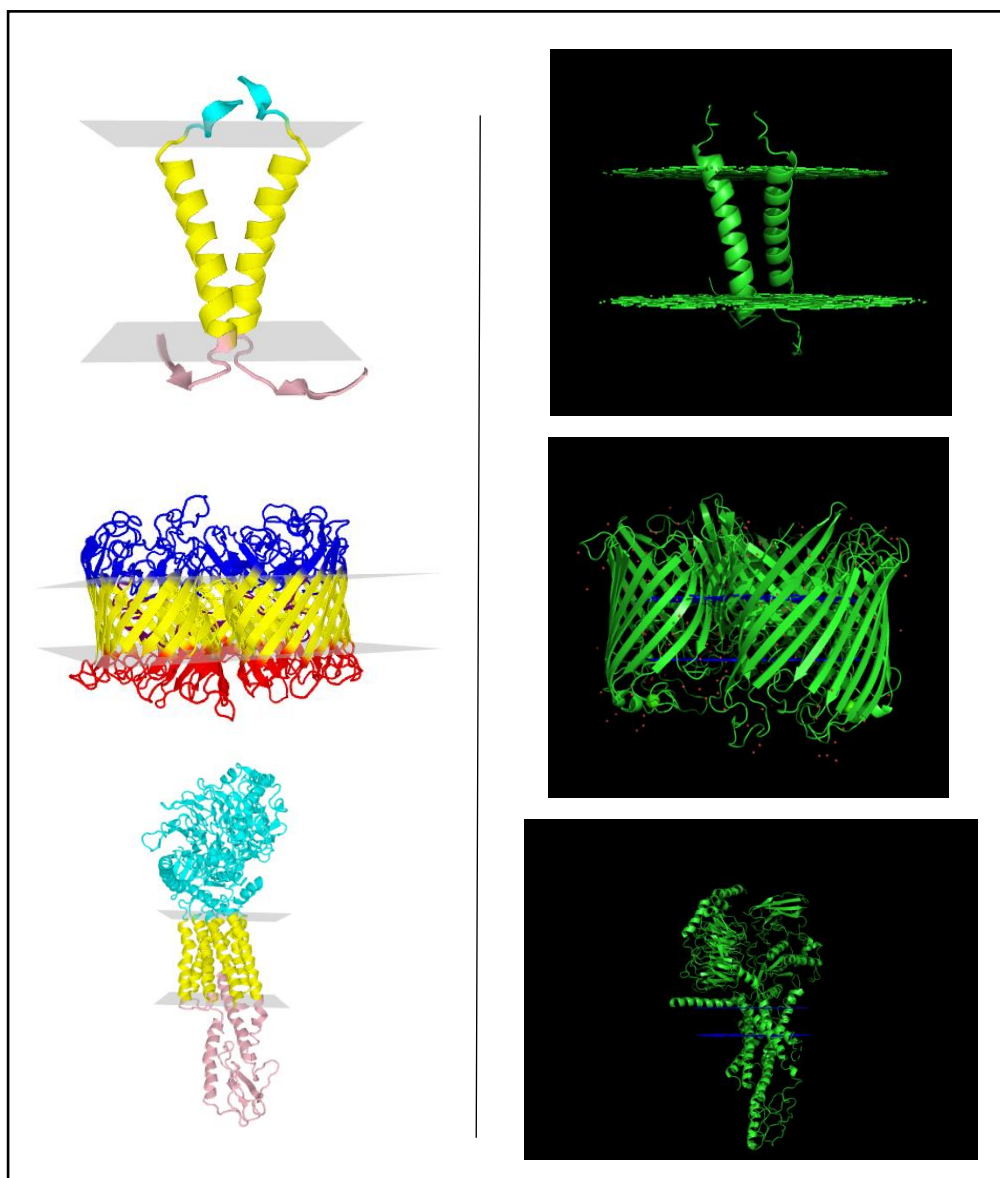


Figure. Localisation de la membrane dans la base de données PDB_TM (gauche) comparé à la localisation obtenue avec l'outil développé pendant ce projet (droite). ID des protéines de haut en bas : 2n90, 1a0s et 6whc. Les plans de la membrane sont représentés en bleues pour 1a0s et 6whc sur les images de droite.

L'outil proposé permet dans certains cas de déterminer la position de la bicouche lipidique par rapport à la structure des protéines transmembranaires avec une précision relativement proche de celle indiquée dans la base de données PDB_TM. Cette méthode se révèle performante pour l'étude des protéines dont la résolution est suffisamment haute, sans présence de chaîne discontinue ou de structure complexe tel que des cavités. L'unique utilisation des coordonnées atomiques des Calphas semble être une approche efficace pour obtenir une approximation satisfaisante.

L'implémentation de la fonction objective, relativement flexible, permet également de tester aisément de nouveau système de scoring pouvant être présenté dans la littérature afin de comparer les résultats obtenus.

La représentation de la membrane par l'ajout d'atomes dans le fichier PDB édité est acceptable pour constater sa position par rapport à la protéine, mais nécessite préférentiellement l'utilisation d'un logiciel de visualisation permettant de colorer chacun des plans (PyMOL suggéré).

La rapidité d'exécution du programme peut représenter un frein dans son utilisation, et dépend principalement du nombre de direction quadrillée dans l'espace par la génération des points sur la sphère. Cette durée est également susceptible de s'allonger lorsque la structure des protéines se complexifie.

En optimisant le temps d'exécution de cet outil, celui-ci pourrait représenter une première approche fiable pour prédire la topologie de la bicouche lipidique dans la structure des protéines transmembranaires. Néanmoins, l'algorithme ne permet pas de différencier les protéines membranaires et globulaires, il serait donc nécessaire de savoir au préalable si la protéine comporte des régions transmembranaires.

ANNEXE

```
1 #Argument Parser pour récupérer les arguments entrés par l'utilisateur.
2 argparser = argparse.ArgumentParser(description="Extract chain from a PDB file")
3 argparser.add argument('infile', help="Path to input file (PDB)")
4 argparser.add argument('outfile', help="Path to output file (PDB)")
5 argparser.add argument('--model', type=int, default=1, help=
6 "Modèle de la protéine étudiée")
7 argparser.add argument('--vector nb', type=int, default=30, help=
8 "Nombre de vecteurs échantillonnés aléatoirement")
9 args= argparser.parse args()
10
11 PDBFile=args.infile
12 OutputFile=args.outfile
13 PDBName=os.path.basename(args.infile)
14 PDBid=PDBName.split(".")[0]
15
16 #Lis le fichier PDB et le passe au programme DSSP, contient la valeur relative de la surf
17 ace accessible (ASA) de chaque CA
18 p = PDBParser()
19 structure = p.get structure(PDBid, PDBFile)
20 model protein = structure[args.model-1]
21 dssp = DSSP(model protein, PDBFile, dssp='mkdssp')
22
23 # Créer un dictionnaire DSSP à partir du fichier PDB, contient les Residus IDs de chaque
24 CA.
25 dssp tuple=dssp dict from pdb file(PDBFile)
26 dssp dict = dssp tuple[1]
27 id tuple = [x[1] for x in dssp dict]
28 id=[y[1] for y in id tuple] #Récupère les Résidus IDs dans une liste
29
30 df DSSP = pd.DataFrame(dssp) #Créer un DataFrame à partir du fichier DSSP de départ.
31 df column=[0,1,3]
32 df DSSP = df DSSP.iloc[:,df column]
33 #Garde seulement les IDs par défaut, le nom des acides aminés et les valeurs de relative
34 ASA.
35 df DSSP.columns=['#','amino acid','Relative ASA']
36 df DSSP['#']=id
37 #Remplace les IDs par défaut du fichier DSSP par les Residus IDs stockées dans le diction
38 naire DSSP.
39 #print (df.loc[[6]])
40 print(df DSSP)
41
42 threshold = 0.20
43 #Seuil minimum de la relative ASA pour considérer qu'un résidu est accessible au solvant.
44 residue hydrophobic=['F','G','I','L','M','V','W','A','C','S','H']
45 #Classification binaire des résidus selon Postic et al., 2015
46 residue hydrophilic=['D','E','K','N','P','Q','R']
47
48 #Conserve seulement les CA si leur résidu est accessible au solvant et classe les hydroph
49 obes et les hydrophiles dans deux listes séparées
50 accessible ca = df DSSP.loc[(df DSSP['Relative ASA']>threshold)]
51 accessible ca hydrophobic = accessible ca.loc[accessible ca['amino acid'].str.contains(
52 '|'.join(residue hydrophobic))]
53 accessible ca hydrophilic = accessible ca.loc[accessible ca['amino acid'].str.contains(
54 '|'.join(residue hydrophilic))]
55
56 #Récupère les résidus IDs des CA accessibles au solvant pour les résidus hydrophobes et h
57 ydrophiles.
58 accessible residue hydrophobic= accessible ca hydrophobic['#'].tolist()
59 accessible residue hydrophilic= accessible ca hydrophilic['#'].tolist()
60
61 #Lecture du fichier PDB en entrée pour récupérer les coordonnées de tous les CA accessibl
62 es au solvant.
63 #Ces coordonnées ne sont pas accessibles depuis le fichier DSSP ou le dictionnaire DSSP.
64 io = PDB.PDBIO
65 struct = p.get structure(PDBid,PDBFile)
66 model = structure[args.model-1]
67 #Ne parcourt que le modèle de la protéine souhaité, par défaut le 1er modèle.
68
69 #Liste les coordonnées x,y,z des CA des résidus hydrophobes et hydrophiles.
70 ca hydrophobic coord list = []
71 ca hydrophilic coord list = []
72
73 for chain in model:
74     for residue in chain:
75         res id=residue.get full id()[3][1]
76         #On ne récupère que le ID du résidu parmi tous les IDs obtenus avec le get_full id
77         if res id in accessible residue hydrophobic :
78             ca=residue["CA"]
79             ca hydrophobic coord=ca.get coord()
80             #Récupère uniquement les coordonnées des CA.
81             ca hydrophobic coord list.append(ca hydrophobic coord)
82
83             elif res id in accessible residue hydrophilic :
84                 ca=residue["CA"]
85                 ca hydrophilic coord=ca.get coord()
86                 ca hydrophilic coord list.append(ca hydrophilic coord)
87
88 #Calcul le centre de masse de la protéine avec les coordonnées de tous les CA des résidus
89 accessibles au solvant.
90 ca coord list = np.array(ca hydrophilic coord list + ca hydrophobic coord list)
91 com=center of mass(ca coord list)
92
93 #Transforme les coordonnées de tous les CA par rapport aux coordonnées du centre de masse
94 (qui représentera le point (0,0,0))
95 transformed coord hydrophilic=transform coordinates(ca hydrophilic coord list,com)
96 transformed coord hydrophobic=transform coordinates(ca hydrophobic coord list,com)
97
98 #Détermine les vecteurs (a,b,c) à partir de l'échantillonnage de points sur sphère.
99 vecteurs=sample spherical(args.vector nb,3)
100 #Création d'un plan orthogonal au vecteur passant par le point (0,0,0).
101 orthogonal planes=get orthogonal planes(vecteurs,x,y,z,d)
```

```

1 optimal_position = []
2 optimal_score_direction = []
3
4 #Boucle "for" pour trouver le meilleur score obtenu dans chaque direction, et à quelle p
osition il a été obtenu.
5 for i in orthogonal_planes :
6     list_score=[]
7     position=0 #The initial position is always 0
8
9     #Donne les équations des deux plans de la membrane à la position 0.
10    membrane1, membrane2 = get_membrane_planes(i)
11
12
13    #Calcul le score obtenu avec la membrane à la position 0.
14    score, nb_hydrophobic_atom_inside=objective_function(membrane1, membrane2,
transformed_coord_hydrophilic, transformed_coord_hydrophobic)
15    list_score.append(score)
16    #print (atoms inside)
17
18
19    #Déplace la membrane dans la direction du vecteur par tranche de 2 Å tant que le nombre
de CA hydrophobes entre les deux plans > 0.
20    while nb_hydrophobic_atom_inside > 0:
21        position += 2
22
23        #Donne les équations des plans de la membrane à la position donnée.
24        membrane1, membrane2 = get_membrane_planes(i,pos=position)
25
26        #Calcul le score obtenu avec cette membrane à la position donnée.
27        score, nb_hydrophobic_atom_inside=objective_function(membrane1, membrane2,
transformed_coord_hydrophilic, transformed_coord_hydrophobic)
28        list_score.append(score)
29
30    #Trouve le score maximum obtenu pour une membrane dans cette direction, et à quel index
de la liste il se trouve.
31    max_score_direction=max(list_score)
32    idx_max_score = list_score.index(max_score_direction)
33
34    #Récupère le score maximal et l'index associé pour chaque direction.
35    optimal_position.append(idx_max_score)
36    optimal_score_direction.append(max_score_direction)
37
38    #print (optimal position)
39
40    #Liste les équations des plans optimales après épaississement dans chaque direction, et
leur score associé.
41    optimal_membrane1_list=[]
42    optimal_membrane2_list=[]
43    optimal_score=[]
44
45    #Boucle "for" pour trouver le meilleur score obtenu après épaississement à la position o
ptimale de chaque direction.
46    for plan,idx_position,score in zip(orthogonal_planes,optimal_position,
optimal_score_direction) :
47        #for each membrane with the highest number of atoms, we increase their thickness 1 Å by
1 Å as long as the score increase
48        add_to_pos=2 #correspond to the size of the pas during the sliding of the membrane
49        thickness = 14
50        max_score=score
51
52        list_mb1=[]
53        list_mb2=[]
54        list_score=[]
55
56        thickness +=1 #Incrémenter l'épaisseur de la membrane de 1 Å par 1 Å.
57        position = 0 + add_to_pos*idx_position
58
59        #Retrouve la position à laquelle la meilleure membrane a été trouvée dans cette directio
n.
60
61        #Calcul les nouvelles équations et le score des plans de la meilleure position après augm
entation de son épaisseur.
62        membrane1, membrane2 = get_membrane_planes(plan, thickness= thickness,pos= position)
63        score,nb_hydrophobic_atom_inside=objective_function(membrane1, membrane2,
transformed_coord_hydrophilic, transformed_coord_hydrophobic)
64
65        list_mb1.append(membrane1)
66        list_mb2.append(membrane2)
67        list_score.append(score)
68
69
70        #Si le score de la membrane après épaississement est supérieur au score précédent, conti
nue d'épaissir la membrane tant que c'est le cas.
71        if score>=max_score:
72            while score >= max_score :
73                max_score=score
74                thickness += 1
75                membrane1, membrane2 = get_membrane_planes(plan, thickness=thickness,pos=
position)
76                new_score,nb_hydrophobic_atom_inside=objective_function(membrane1, membrane2
, transformed_coord_hydrophilic, transformed_coord_hydrophobic)
77
78                list_mb1.append(membrane1)
79                list_mb2.append(membrane2)
80                list_score.append(new_score)
81
82                if new_score >= score :
83                    max_score = score
84                    score = new_score
85
86        #Si le nouveau score est inférieur au précédent, récupère les équations et le score de l
a membrane précédente.
87        else :
88            membrane1, membrane2 = get_membrane_planes(plan, thickness=thickness,pos
= position)
89            #if the score doesn't increase for the new membrane, get the equation of the previous mem
brane + its score
90            optimal_membrane1_list.append(list_mb1[-2])
91            optimal_membrane2_list.append(list_mb2[-2])
92            optimal_score.append(list_score[-2])
93            break
94
95        #Si le premier épaississement n'a pas amélioré le score, garde les équations et le score
obtenus sans épaississement.
96        else :
97            membrane1, membrane2 = get_membrane_planes(plan,pos= position)
98            #keep the original thickness if the incrementation doesn't increase the score from the begi
nning
99            score, nb_hydrophobic_atom_inside=objective_function(membrane1, membrane2,
transformed_coord_hydrophilic, transformed_coord_hydrophobic)
100            optimal_membrane1_list.append(membrane1)
101            optimal_membrane2_list.append(membrane2)
102            optimal_score.append(score)
103
104        #Trouve le score maximal et l'index associé obtenus après épaississement de toutes les m
embranes à la meilleure position de chaque direction.
105        max_score_membrane=max(optimal_score)
106        idx_score_optimal = optimal_score.index(max_score_membrane)
107
108        #Retourne les équations des plans de la membrane optimale pour la protéine.
109        optimal_membrane1=optimal_membrane1_list[idx_score_optimal]
110        optimal_membrane2=optimal_membrane2_list[idx_score_optimal]
111        print (optimal_membrane1, optimal_membrane2)

```



```

1 #Récupération de chaque argument au sein des équations.
2 #Doit utiliser findall car les équations sont composées de symbole sympy, l'ordre des arg
3 uents retournés par .arg[] change entre les équations.
4 str_mb1_equation=str(optimal_membrane1)
5 str_mb2_equation=str(optimal_membrane2)
6 str_mb1_parameters= reg.findall(r"(?:\+|s*|\-|s*)?[-+]?d*\.*d+", str_mb1_equation)
7 str_mb2_parameters= reg.findall(r"(?:\+|s*|\-|s*)?[-+]?d*\.*d+", str_mb2_equation)
8 mb1_parameters=[float(str(i).replace(' ','')) for i in str_mb1_parameters]
9 mb2_parameters=[float(str(i).replace(' ','')) for i in str_mb2_parameters]
10
11 a,b,c=mb1_parameters[0],mb1_parameters[1],mb1_parameters[2]
12 d1=mb1_parameters[3]
13 d2=mb2_parameters[3]
14
15 normal_vec=(a,b,c)
16
17 #Génération de points (x,y,z) random situés dans le plan de la 1ere couche de la membran
18 e.
19 r=d1/c
20 refpoint=(0.0,0.0,r)
21 plan1=Plane(normal_vec,d1,refpoint,30)
22 random_plane_points = plan1.create_random_points(1000)
23 #https://pypi.org/project/random-geometry-points/
24
25 #Génération de points random (x,y,z) situés dans le plan de la 2ere couche de la membran
26 e.
27 r2=d2/c
28 refpoint2=(0.0,0.0,r2)
29 plan2=Plane(normal_vec,d2,refpoint2,30)
30 random_plane_points_2 = plan2.create_random_points(1000)
31
32 #Lecture du fichier PDB en entrée puis transformation des coordonnées de chaque atome par
33 rapport au centre de masse.
34 parser = PDB.PDBParser()
35 io = PDB.PDBIO()
36 struct_transform = parser.get_structure(PDBid,PDBFile)
37 model_transform=struct_transform[args.model -1]
38 rotation_matrix = PDB.rotmat(PDB.Vector([0, 0, 0]), PDB.Vector([0, 0, 0]))
39 #Matrice neutre, pas de rotation des coordonnées souhaitée.
40 transform=(-com[0],-com[1], -com[2])
41
42 for chain in model_transform:
43     for residue in chain:
44         for atom in residue:
45             atom.transform(rotation_matrix, transform)
46
47 #Retourne un fichier PDB avec les coordonnées de tous les atomes transformées.
48 io.set_structure(model_transform)
49 PDBTransform= os.path.join(OutputFile,'{0}_model{1}_transform.pdb'.format(PDBid, args
50 .model))
51 io.save(PDBTransform)
52
53 #Lis le fichier contenant les nouvelles coordonnées avec le module Biotite.
54 atom_array = strucio.load_structure(PDBTransform)
55
56 #Pour chaque point généré aléatoirement sur la 1ere couche de la membrane , créé un hétéro
57 oatome possédant ses coordonnées.
58 for i in random_plane_points :
59     atom = struc.Atom(
60         coord = [i[0],i[1],i[2]],
61         chain_id = "N",
62         res_id = atom_array.res_id[-1] + 1,
63         #Le résidu id de l'atome est égale au dernier résidu ID du fichier +1
64         res_name = "MB1",
65         #Chaque atome créée pour la 1ère couche possède le nom de résidu MB1.
66         hetero = True,
67         atom_name = "CA",
68         element = "C"
69     )
70     atom_array += struc.array([atom])
71 #Ajout de l'hétéroatome à la suite des atomes présent dans le fichier PDB.
72
73 #Pour chaque point généré aléatoirement sur la 2eme couche de la membrane, créé un hétéro
74 ome possédant ses coordonnées.
75 for j in random_plane_points_2 :
76     atom = struc.Atom(
77         coord = [j[0],j[1],j[2]],
78         chain_id = "N",
79         res_id = atom_array.res_id[-1] + 1,
80         res_name = "MB2",
81         #Chaque atome créée pour la 1ère couche possède le nom de résidu MB2.
82         hetero = True,
83         atom_name = "CA",
84         element = "C"
85     )
86     atom_array += struc.array([atom])
87 #Ajout de l'hétéroatome à la suite des atomes présent dans le fichier PDB.
88
89 #Retourne un fichier PDB édité avec les hétéroatomes représentant chaque couche de la mem
90 brane ajoutés à la structure de la protéine.
91 PDBEdited=os.path.join(OutputFile,'{0}_model{1}_edited.pdb'.format(PDBid,args.model))
92 strucio.save_structure(PDBEdited, atom_array)

```