

Internship project 2015



# Computational Protein Design of a PDZ Domain using Proteus

In the Biochemistry laboratory of the Ecole Polytechnique

From January 12 to June 30, 2015

By Xingyu CHEN

Master2 student (2014/2015) of In Silico Drug Design, University Paris Diderot – Paris7  
Add.: 166 Avenue de La Division Leclerc, 92160 Antony, France  
Phone: +33 (0)6 03 81 15 88  
Email: xingyu809@gmail.com

Supervised by Professor Thomas SIMONSON  
Add.: Department of Biology, Laboratoire de Biochimie (UMR 7654 du CNRS), Ecole Polytechnique, 91128 Palaiseau, France  
Phone: +33 (0)1 69 33 48 60 (fax: -- 49 09)  
Email: thomas.simonson@polytechnique.fr

## **ACKNOWLEDGMENTS**

I thanks to all members in BIOC and Bioinformatics team members, especially my supervisor Thomas Simonson, who helps me patiently and suggests on my work. I gratefully acknowledge David Mignon for his help in solving technical problems. I gratefully acknowledge Nicolas Panel and Karen Druart for their discussions in my project. I gratefully acknowledge Anke Steinmetz, who always encourages me and gives me her expert suggests in my project.

## TABLE OF CONTENTS

ACKNOWLEDGMENT -----	i
TABLE OF CONTENTS -----	ii
1. INTRODUCTION -----	1
1.1 Context -----	1
1.2 PDZ Structure and Classification -----	1
1.3 Experimental sequences of PDZ peptide ligand -----	1
1.4 Internship subject -----	2
2. MATERIALS & METHODS -----	2
2.1 General CPD method -----	2
2.2 Application systems -----	3
2.3 Reference energies, method 1: syndecan peptide as the unfolded model -----	5
2.4 Reference energies, method 2: empirical model -----	5
2.5 Simulations of the Tiam1 PDZ:peptide complex -----	7
3. RESULTS -----	8
3.1 Sequence alignments and experimental amino acid compositions -----	8
3.2 Simulations of the syndecan peptide -----	8
3.3 Simulations of the Tiam1: peptide complex – version1 -----	8
3.4 Optimization of empirical reference energies -----	9
3.5 Simulations of the Tiam1: peptide complex – version2 -----	10
4. DISCUSSION -----	11
5. REFERENCES -----	12
6. APPENDIX -----	32

# 1. INTRODUCTION

## 1.1 Context

The Rho (Ras homologous) GTPase (Guanosine Triphosphatase) family is well known for its regulation of actin cytoskeleton dynamics and is crucial for cell migration. Rho GTPase misregulation contributes to cancer cell invasion, tumorigenesis and metastasis<sup>[1]</sup>. The T-cell lymphoma invasion and metastasis gene 1 (Tiam1), a multiple domain protein, is identified as a guanine nucleotide exchange factor (GEF) that activates the Rho-family GTPase Rac1 (Ras-related C3 Botulinum Toxin Substrate 1)<sup>[2]</sup>. Tiam1 has important roles in cell biological function, including cell polarity, integrity of adherent junctions<sup>[3]</sup>, tight junctions<sup>[4-6]</sup>, and cell-matrix interactions<sup>[7, 8]</sup>, whose deregulation has been documented in many cancers and Tiam1/Rac1 signalling has been implicated in the oncogenic transformation of cells. Tiam1 can be inhibited by small molecules that bind to its PDZ domain, and block or down regulate the interaction with its target proteins. We focus on the PDZ domain of the Tiam1 protein for searching peptide inhibitors.

## 1.2 PDZ Structure and Classification

The common structure of PDZ domains comprises six  $\beta$  strands and two  $\alpha$  helices. A groove between the  $\beta$  strand and the  $\alpha$  helix can accept a peptide ligand at the C-terminal binding as a beta sheet extend (*figure1*). The PDZ domains are classified into 3 classes according to the sequences of their binding peptide ligand. All 3 types of PDZ domain prefer a ligand with a hydrophobic residue at 0 position from C-terminal. At the -2 position of the ligand, the class II PDZ domain prefers a hydrophobic residue, whereas the class I PDZ domain prefers the residue S and T, and the class III PDZ domain prefers the residue D and E. The PDZ domain of Tiam1 is the class II PDZ domain. Some precedent studies show that, ligands of Tiam1 PDZ domain has a consensus amino acid sequence motif of  $[Y/F/G]_2-[Y/F/H/W]_1-[W/F/A]_0-COOH^{[9, 10]}$ . In this project, we use our computational protein design (CPD) method to generate mutants of ligands at the positions -4 to 0, in Tiam1 PDZ complexes.

## 1.3 Experimental sequences of PDZ peptide ligand

There are not many experimental structures of class II PDZ. We collected two experimental sequence sets of PDZ peptide ligand. We searched the crystallography structures of PDZ domain in Protein Data Bank, and found 143 PDZ structures that are formed in complex. Among these 143

structures, there are only 30 structures of class II PDZ (the PDB codes and the ligand sequences are presented in *table1*). The first set collect the 26 sequences of ligand peptide from class II PDZ complexes. The second set is from a Tiam1-binding peptide library, which contains 69 peptide sequences in total<sup>[11]</sup>. The sequences in library are “YAA<sub>-4</sub>X<sub>-3</sub>X<sub>-2</sub>X<sub>-1</sub>X<sub>0</sub>”, in which X represents a given residue type. The sequence logos of the ligand at position -4 to 0 are presented in *figure2*. The size of residue logo represents its population at each position.

#### **1.4 Internship subject**

The global goal of the PDZ project is design the peptide inhibitors of Tiam1 PDZ domain for anticancer applications. In my internship project, we want to evaluate the performance of the CPD method, which is developed in our group, by mutagenesis of ligand in Tiam1 PDZ complexes. To this end, we developed different protocols to optimize the parameters of CPD method, the reference energies. We applied these optimized parameters and compared the calculated sequences with experimental sequences.

## **2. MATERIALS & METHODS**

### **2.1 General CPD method**

To perform CPD, our laboratory developed the Proteus software package, which is implemented by xplor, shell, and perl scripts. The method is already published<sup>[12]</sup>. There are three main steps: system construction, energy matrix calculation, and sequence exploration.

#### **2.1.1 Structural model**

The folded protein or peptide is modeled using a fixed backbone and a discrete rotamer library for the side chains; solvent is modelled implicitly. For the unfolded state, we use a extended peptide model, for examples, ALA-X-ALA. We consider that each residue X in the protein interacts only with the nearby backbone and solvent, not with other side chains.

Each residue in the protein is defined as one of three types: residues with fixed side chain are called as “frozen”; the residues with flexible side chains that cannot mutate are called “inactive”, and the residues that can be mutated are called as “active”.

#### **2.1.2 Energy function**

The energy function is calculated using an MM/GBSA model. We approximate the energy further by decomposing it into residue terms. The total energy equals the sum of self-interaction energy of

residue i ( $E_{ii}$ ) plus the sum of pairwise interaction energy between residue i and residue j ( $E_{ij}$ ).

$$E = \sum E_{ii} + \sum E_{ij} \quad (1)$$

For the unfolded state, we introduce “reference energies” for each residue type X ( $E_X$ ), which represent the contribution of amino acid X in the unfolded state. Each value includes a contribution  $E_X^{\text{pept}}$  from an extended peptide model plus an empirical contribution  $e_x$  that can be adjusted. The unfolded energy of protein ( $E_{\text{unfolded}}$ ) is the sum of the reference energies of each position i (equation (2)).

$$E_{\text{unfolded}} = \sum E_X(i) \quad (2)$$

### 2.1.3 Energy matrix

The interaction energies are stored in an energy matrix, used as a look-up table. The self-interaction energy term ( $E_{ii}$ ) and the reference energy ( $E_X$ ) are saved in the diagonal of the matrix; the pairwise interaction energy terms ( $E_{ij}$ ) are saved in off diagonal terms.

### 2.1.4 Mutation space exploration – MC

The mutation space exploration is achieved by proteus program. In this project, we used a Monte Carlo (MC) method. MC generates sequences and conformations according to a Boltzmann distribution that depends on the folding free energy:

$$E_{\text{folding}} = E_{\text{folded}} - E_{\text{unfolded}} \quad (3)$$

For selected positions of the protein (or peptide) sequence, we can do single point mutations or multiple mutations and search the rotamer conformations, or only search the rotamer conformations. The results will be analysed for testing and improving the reference energies  $E_X$ , or for selecting the interesting mutants.

## 2.2 Application systems

### 2.2.1 The syndecan peptide

For preparing the system, we used the Tiam1:peptide Xray complex (PDB: 4GVD\_A\_D) and minimized slightly with the procedure in below:

1. 200 steps minimization of Steepest Descent (SD) method with 50 kcal protein constraints.
2. 200 steps minimization of SD method with 10 kcal protein constraints.
3. 200 steps minimization of SD method and of Adopted Basis Newton-Raphson) method respective, and with 2 kcal protein constraints.
4. 200 steps minimization of sd method and abnr method respective, and with released protein.

The histidines in complex are treated by the software “PROPKA” and “Reduce”. The two histidines in protein are protonated at delta position with the pka value equals to 6.18 for histidine at position 844 and 6.34 for the one at position 847.

We firstly worked with the extended Syndecan1 peptide. The peptide structure was extracted from the minimized complex structure. The sequence is “TKQE<sub>4</sub>E<sub>-3</sub>F<sub>-2</sub>Y<sub>-1</sub>A<sub>0</sub>-COOH”. For the energy matrix calculation, we used solute a dielectric constant of either 4 or 8 for the GB energy term. The atomic surface coefficients were (in kcal/mol/Å<sup>2</sup>): alkane atoms = -0.005; hydrophobic atoms = 0; polar atoms = -0.008; aromatic atoms = -0.012; ionic atoms = -0.009. Each diagonal ( $E_{ii}$ ) and off-diagonal ( $E_{ij}$ ) term in the energy matrix was computed after 15 steps of Powell minimization (per pair). We ran MC simulations with 5 million steps, to explore single point mutations with all 20 amino acid types allowed at the last five peptide positions. The MC simulation uses a “mono walker” with the temperature equals to 0.6K. The movement probabilities of rotamer (*Rot\_Proba* and *Rot\_Rot\_Proba*) are equal to 0.5. The movement probabilities of mutation (*Mut\_Proba*, *Mut\_Mut\_Proba* and *Mut\_Rot\_Proba*) are equal to 0.1.

### **2.2.2 Tiam1 PDZ protein**

The second system is the Tiam1 PDZ protein. The structure was extracted from Tiam1:peptide Xray complex (PDB: 4GVD\_A\_D) after the minimizations. For the energy matrix, we used the same parameters as for Syndecan1 above. We ran MC simulations with 10 million steps, to explore multiple mutations over all the positions of protein (in the absence of the peptide ligand). For the MC simulation, we used “8 walkers” with the temperatures (unit: K) equal to 3, 2, 1.333, 0.888, 0.592, 0.395, 0.263, and 0.175. The movement probabilities of rotamer and mutation are the same as in above. We only selected the sequences generated with the temperature 0.592.

### **2.2.3 Cask PDZ protein**

Our third system is built with the PDZ protein of hCask from the Xray structure (PDB code: 1KWA\_A). The missing residues in the file pdb are added by “scwrl”. The histidines are treated as protonated at both delta and epsilon positions. We used the same procedure and parameters as for Tiam1 PDZ protein.

### **2.2.4 Sequence alignment and amino acid composition of PDZ proteins**

To calculate the experimental amino acid composition of our proteins, we collected two sets of natural sequences from UNIPROT: one with Tiam1 homologues and the other with hCask homologues. The homologues with sequence identities compared to the target protein of between

60% and 85% were selected. There were 50 homologues of Tiam1 and 126 of Cask.

### 2.3 Reference energies, method 1: syndecan peptide as the unfolded model

A first method to calculate the reference energies relies on the extended Syndecan1 peptide as the unfolded state backbone model. From our MC peptide simulations (above), we used the reference energies, which are obtained from the precedent study with a peptide fragment in complex with MHC class II protein ( $\{E_X^i\}_{MHC}^{114}$ ). As MC method uses folding free energy to generate mutants (equation (2) and (3) above), and as our peptide is extended, the folded energy ( $E_{folded}$ ) in our case represents the extended peptide energy of mutants ( $E_{mut}$ ). The energies at unfolded state are represented by the sum of initial reference energies ( $\sum_i E_{Xini}^i$ ). Thus, the MC energy can be represented in equation (4). We calculated the energies of all point mutant respect to alanine (A) with the equation (5) ( $E_{mut}^{Xi}$  represent the extended energy of the mutant with given residue type X at position i,  $E_{mut}^{Ai}$  represent the extended energy of the mutant with residue type A at position i).

The equation (5) can be transformed in equation (6).

$$E_{MC} = E_{mut} - (\sum_{j \neq i} \{E_X^j\}_{MHC} + \{E_X^i\}_{MHC}) \quad (4)$$

$$\{\delta E_{mut}^{pept}\}_X^i = (E_{mut}^{Xi} - \sum_{j \neq i} \{E_X^j\}_{MHC} - \{E_X^i\}_{MHC}) - (E_{mut}^{Ai} - \sum_{j \neq i} \{E_A^j\}_{MHC} - \{E_A^i\}_{MHC}) \quad (5)$$

$$\{\delta E_{mut}^{pept}\}_X^i = E_{mut}^{Xi} - E_{mut}^{Ai} - \{E_X^i\}_{MHC} + \{E_A^i\}_{MHC} \quad (6)$$

We obtained energies for all possible single point variants at the last five peptide positions. For each sidechain type X, the energies of the point mutant difference relative to Ala and their averaged value (at the same peptide position) was then taken to represent  $E_X$ . For Ala, we can arbitrarily set  $E_X = 0$ .

### 2.4 Reference energies, method 2: empirical model

#### 2.4.1. General method

Here, we use a simple extended peptide model to descript the unfolded state. But this simple model needs to be supplemented by an empirical correction. Thus, we add an empirical energy correction  $e_X$  to each  $E_X$ , to avoid the too abundant occurrence of a given residue type X:

$$E_X = E_X^{ini} + e_X \quad (7)$$

$e_X$  depends only on the amino acid type. We iteratively modify  $e_X$  using the Boltzmann-like relation<sup>[12, 14]</sup>:

$$e_X^{new} = e_X^{old} + 0.5 \ln f_X^{exp} / f_X^{cal} \quad (8)$$

where  $f_X^{exp}$  represents the experimental composition of amino acid type X (obtained by counting the

occurrences of X in a set of homologs, see above).  $f_X^{\text{cal}}$  represents the theoretical composition of amino acid type X, obtained from one or more MC simulations. We apply (8) repeatedly; the iterations will be stopped when  $f_X^{\text{cal}}$  of all amino acid types have converged close to  $f_X^{\text{exp}}$ . A flowchart of the procedure is presented in *figure3*. In general case, we have 20 acid amino types and  $e_X$  values to adjust. But to decrease the number of residue types, we can also group similar residues, so that they share a single  $e_X$  value (but different  $E_X^0$  values). To refine the model, we can assign different  $e_X$  values to positions that are buried or exposed (in the folded state).

#### **2.4.2 Implementation details for Tiam1, version 1**

The experimental amino acid compositions are computed with the sequences of the 50 Tiam1 homologues. The calculated amino acid compositions are generated by a 500 million step MC where all protein positions mutate freely. The MC simulation uses a “multi-walker”, Replica exchange protocol, with 8 walkers (the 8 temperatures as in section 2.2.2). For the initial  $E_X$  values, we used the lowest self-interaction energies ( $E_{ii}$ ) of each residue type, averaged over all the protein positions. We set the initial  $e_X$  values to zero. We iteratively updated the reference energies until convergence (10 iterations for eps8, and 15 iterations for eps4 but not yet converged).

#### **2.4.3 Implementation details, version 2**

The second implementation is with Tiam1 and hCask protein together. We classified residues into 7 groups according to the matrix BLOSUM62<sup>[13]</sup> (*table2*).

Each group shares a single  $e_X$  value (zero, initially). We also separated the buried and exposed positions using a surface burial fraction threshold equals to 0.2, allowing them to have different  $e_X$  values. There are 33 buried and 51 exposed positions in Tiam1, and 34 buried and 40 exposed positions in Cask. Thus, each residue type has two reference energy values: buried reference energy ( $E_{Xb}$ ) and exposed reference energy ( $E_{Xe}$ ). The iteration formulas are:

$$\text{Buried residues: } E_{Xb}^{\text{new}} = E_{Xb}^{\text{old}} + 0.5 \ln(f_{gb}^{\text{exp}}/f_{gb}^{\text{cal}})$$

$$\text{Exposed residues: } E_{Xe}^{\text{new}} = E_{Xe}^{\text{old}} + 0.5 \ln(f_{ge}^{\text{exp}}/f_{ge}^{\text{cal}})$$

For the experimental compositions of amino acid groups ( $f_g$ ), we calculated separately for the 50 homologues of Tiam1 PDZ and the 126 homologues of Cask PDZ, and averaged the two. The calculated compositions of amino acid groups are computed from 10 separated MC simulations: 5 for Tiam1 and 5 for Cask, at each iteration. For the 5 MC simulations with each protein, we assigned in order 1 per 5 positions as active to each simulation for the multiple mutation explorations, and the rest 4/5 inactive positions with the side chain flexible. The MC simulations

used the “8 walkers” as in section 2.4.2 above. The initial  $E_x$  are used the delta mutant energies, which averaged over the 5 active positions of peptide from the section 2.3. The shell script for reference energy optimization is in *Appendix*.

## 2.5 Simulations of the Tiam1 PDZ:peptide complex

To test the quality of the reference energies, which are obtained from the different methods above, we applied these reference energies to explore ligand mutants with Tiam1 PDZ complexes. For the simulations, we used 3 different starting structures:

1. The Xray structure: Tiam1 PDZ complex binding with Syndecan1 (PDB: 4GVD\_A\_D) and we minimized slightly with the procedure describing in the section 2.2.1.
2. The Tiam1 PDZ complex with the peptide ligand, which was selected from the Tiam1-binding library<sup>[11]</sup> with the best binding affinity score. The sequence of peptide is “YAAE<sub>-4</sub>K<sub>-3</sub>Y<sub>-2</sub>W<sub>-1</sub>A<sub>0</sub>”. The structure is obtained by 20ns molecular dynamics simulation.
3. Tiam1 PDZ complex binding the syndecan1 mutant with Ala mute to Phe at position 0. We used a 100ns molecular dynamics simulation. The structure used is the structure most similar than the average structure from the molecular dynamics simulation.

For each structure, we built two systems both with the protein flexible (side chain flexible) and frozen (side chain fixed), and calculated the energy matrixes. In each system, ligand peptides are always modeled with the first 3 positions inactive (side chain flexible), and the last 5 positions active (will be mutated). The dielectric constant was either 4 or 8 (the other parameters are as the parameters used in the simulations above).

For sequence exploration, we used the MC method for both multiple mutations and single point mutations at positions -4 to 0 of peptide ligand in Tiam1 PDZ complex. The multiple mutation explorations are achieved by a 10-million-step MC over all active positions. The single point mutations are achieved by 5-million-step MC at each active position. The MC simulations are used the “8 walkers” as described in above.

The interested structures were reconstructed, and minimized by 200 steps using xplor.

### 3. RESULTS

#### 3.1 Sequence alignments and experimental amino acid compositions

The two sets of homologues were aligned separately by the *Clustal Omega* program in *Uniprot* and viewed with *Jalview*. The alignment with 50 Tiam1 homologues is presented in *figure5*. The alignment with 126 Cask homologues is presented in *figure6*. The identify percentage between Tiam1 and Cask is 23.68%. We eliminated the redundant sequences, and the aligned sequences are ordered by their identify percentages. The amino acid compositions of these two homologue sets are represented in *table4*. Results are shown separately for the buried and exposed protein positions and for groups of similar amino acids.

#### 3.2 Simulations of the syndecan peptide

We simulated the syndecan peptide in its extended conformations with the last 5 positions allowed to mutate, either individually or together. For the single point mutations, we obtained all amino acid types at each active position of peptide. The results are presented by sequence logos (*figure4*). The size of the letters corresponds to the residue type population. The types at each mutated position are sorted by their populations. The letter “h” represents the histidine protonated at delta position (Hid) and “j” represents the histidine protonated at epsilon position (Hie). The energies of all point mutants, relative to Ala, and their averaged values are shown in blue bloc of *table3*. These energies are similar at all peptide positions, consistent with its extended structure. They will be applied to Tiam1:peptide complex as the reference energies for the peptide mutation space exploration.

#### 3.3 Simulations of the Tiam1: peptide complex – version1

Next, we considered the peptide in complex with Tiam1. If we assume the extended syndecan peptide is a good model of an unfolded protein, we can use the mutation energies above as reference energies  $E_X$ . We explored mutations in the same five peptide positions, individually or together. We used energy matrices corresponding to three slightly different 3D structures. The results are presented by sequence logos. For the results of multiple mutation explorations, the sequences that have the occurrence upper than 0.2% were aligned. The results are viewed in *Jalview* and show in *figure11-13*.

The first structure we used is the Tiam1:syndecan1 Xray complex (slightly minimized). The logos are presented in *figure7*. The peptide positions -1, -3 and -4 vary extensively, especially

position -1. At position 0 (the C-terminus), we only obtained small residues (Ala, Cys, Ser), not aromatic residues Phe and Trp. At position -2, with a frozen protein, we got many Phe, which is the native residue, and its similar residue His (protonated or not). The results obtained using frozen proteins are more native-like than using flexible proteins, as expected. The MC energies of the single point mutants (using eps8 and frozen protein) at position -2 and 0 are presented in *table6*. The 3D structures of top50-population sequences from multiple mutations with frozen protein and eps8 were reconstructed and minimized. A cross-eye stereo figure of *Pymol* is presented in *figure15* with the structures of top3-population sequences and of Xray.

The second structure is an MD model of Tiam1 bound to a peptide from the Tiam1-binding library<sup>[II]</sup> with the best binding affinity score. The sequence of peptide is “YAAE<sub>4</sub>K<sub>3</sub>Y<sub>2</sub>W<sub>1</sub>A<sub>0</sub>”. The results are presented in *figure8*. The obtained residues at position -1 and -3 are again very diverse, and that of the other 3 active positions are conserved. At position 0, we obtained almost only Ala, except a small population of Ser obtained by multiple mutations with a flexible protein. Many protonated His are found at position -2. At position -4, we obtained a large population of Ala and a little Cys.

Comparing the results with the experimental Tiam1-binding library<sup>[II]</sup> (*figure2b*), we found the results are not very satisfying. The experimental studies have already obtained consensus sequences with Phe at position 0 of Tiam1 PDZ ligands<sup>[II]</sup>. But, with both 3D structures, we only obtained small residues (Ala, Cys, and Ser) at position 0. Thus, we also used a third MD structure, Tiam1 binding the syndecan1 mutant with Ala0 changed to Phe (A0F). The resulting sequence logos are presented in *figure9*. At position 0, we found a lot of His, which has similar structure with Phe, by the simulations with frozen protein, and with flexible protein, the residue Cys has the largest population. We did not regain Phe at position 0. In additional simulations that imposed Phe at position 0, we found the Phe residue type has energy about 2-3 kcal/mol too high to be sampled.

### 3.4 Optimization of empirical reference energies

To improve the peptide design results, we wish to introduce an empirical correction  $e_x$  to the reference energies used above. We explored two different methods to optimize the  $e_x$ ; we only present the more sophisticated one, which gives the best results. With this method, we assign the same  $e_x$  value to a group of similar amino acid types, to reduce the number of parameters. At the same time, to refine the model, we assign different values to positions that are buried or exposed in

the folded structure. We adjust the  $e_x$  iteratively, so that the computed amino acid compositions match the experimental ones. The computed compositions at each iteration are obtained from ten MC simulations, five each of Tiam1 and Cask, with five different sets of 15 active positions.

The optimized reference energies are presented in green bloc of *table3*. The reference energies of buried residues are very similar to those of exposed residues. They are also close to the reference energies obtained by the simulations of the Syndecan peptide. The calculated amino acids compositions, after convergence (15 iterations), are presented in *table5*. The group compositions agree well with the experimental ones, but within each group, some amino acid compositions agree less well with the experimental values.

### 3.5 Simulations of the Tiam1: peptide complex – version2

We next used the improved reference energies to explore mutations of the syndecan peptide bound to Tiam1 (Xray structure). The MC energies of the 20 single point mutants (using eps8 and frozen/flexible protein) at position -2 and 0 are presented in *table6* and *table7*. The sequence log results are presented in *figure10*. We obtained a lot of hydrophobic residues at each position with all simulations. At position 0, we obtained always only small residues. We have got a large population of aromatic groups at positions -2 and -1, and a large population of hydrophobic residues at position -3. These results are consistent with the experimental results.

By the simulation with frozen protein, large populations of aromatic residues are obtained at position -4 too. At position -3, a lot of acidic and polar residues (Glu and Asp) are found, which are consistent with the native residue Glu. The results of single point mutations with both flexible and frozen protein, and the results of multiple mutations with frozen protein, are similar. We can conclude that these results are more robust than those obtained in the previous simulations (no  $e_x$  corrections). For the results of multiple mutation explorations, the sequences that have the occurrence upper than 0.2% were aligned. The results are viewed in *Jalview* and show in *figure14*.

The 3D structures of top50-population sequences from multiple mutations with frozen protein and eps8 were reconstructed and minimized. A cross-eye stereo figure of *Pymol* is presented in *figure15* with the structures of top3-population sequences and of Xray.

We also reconstructed 3D structures of the best 200 conformations of each point mutant (at position 0 and -2) sequence with flexible protein and eps8. The averaged energies over the best 200 conformations and the affinities are presented in *table8*.

## DISCUSSION

We studied Tiam1:peptide complexes using two different models for the unfolded state and two sets of reference energies  $E_x$ . The first model uses the unfolded peptide backbone to represent the unfolded state (of Tiam1). The mutation energies at each peptide position are similar; the position-average values are used for the  $E_x$ . The second method adds an empirical correction. We have two implementations with this method. The first, simpler implementation used Tiam1 (not Cask) and gave results that seem to be of moderate quality (not shown). The second implementation is more sophisticated and the  $e_x$  values are more constrained (use of groups, fewer active positions, use of two proteins). It also distinguishes buried and exposed positions (even though  $E_x$  is supposed to represent the unfolded state, where everything is exposed). The buried  $E_x$  and the exposed  $E_x$  of a given residue type are very similar, and similar to the unfolded peptide model (the  $e_x$  are small) which is encouraging. We can conclude that these reference energies are fairly general, not depending much on the positions or protein.

In the simulations of the Tiam1:peptide complexes, we observe that with a slight difference of the reference energies, the obtained sequence are significantly different. This reflects the exponential dependence of the populations on the  $E_x$  values. With the latest values, the MC results show a good agreement with the experimental peptide sequences and are more robust than the precedent results. But the multiple mutation results with a flexible protein in all the systems are less good, possibly because we did not include the Xray orientation of the sidechains among the possible rotamers. In the next work, we will add these “native” rotamer in our method for the mutation exploration.

## REFERENCES

- [1] Li, H., Peyrillier, K., Kilic, G. & Brakebusch, C. Rho GTPases and cancer. *Biofactors* (2013).
- [2] Habets GG, Scholtes EH, Zuydgeest D, van der Kammen RA, Stam JC, Berns A, Collard JG. Identification of an invasion-inducing gene, Tiam-1, that encodes a protein with homology to GDP-GTP exchangers for Rho- like proteins. *Cell*, 1994, 77(4): 537–549
- [3] Malliri, A., van Es, S., Huvemeers, S. & Collard, J.G. The Rac exchange factor Tiam1 is required for the establishment and maintenance of cadherin-based adhesions. *J Biol Chem* 279, 30092-8 (2004).
- [4] Chen, X. & Macara, I.G. Par-3 controls tight junction assembly through the Rac exchange factor Tiam1. *Nat Cell Biol* 7, 262-9 (2005).
- [5] Mertens, A.E., Rygiel, TP., Olivo, C., van der Kammen, R. & Collard, J.G. The Rac activator Tiam1 controls tight junction biogenesis in keratinocytes through binding to and activation of the Par polarity complex. *J Cell Biol* 170, 1029-37 (2005).
- [6] Nishimura, T. et al. PAR-6-PAR-3 mediates Cdc42-induced Rac activation through the Rac GEFs STEF/Tiam1. *Nat Cell Biol* 7, 270-7 (2005).
- [7] Shepherd, T.R. et al. The Tiam1 PDZ domain couples to Syndecan1 and promotes cell-matrix adhesion. *J Mol Biol* 398, 730-46 (2010).
- [8] Sander, E.E., van Delft, S., ten Klooster, J. P., Reid, T., van der Kammen, R. A., Michiels, F. & Collard, J. G.. Matrix-dependent Tiam1/Rac signaling in epithelial cells promotes either cell-cell adhesion or cell migration and is regulated by phosphatidylinositol 3-kinase. *J. Cell Biol.* 143, 1385–1398. (1998).
- [9] Songyang, Z., Fanning, A. S., Fu, C., Xu, J., Marfatia, S. M., Chishti, A. H., Crompton, A., Chan, A. C., Anderson, J. M., and Cantley, L. C. (1997) Recognition of unique carboxyl-terminal motifs by distinct PDZ domains. *Science* 275, 73–77.
- [10] Tonikian, R., Zhang, Y., Sazinsky, S. L., Currell, B., Yeh, J. H., Reva, B., Held, H. A., Appleton, B. A., Evangelista, M., Wu, Y., Xin, X., Chan, A. C., Seshagiri, S., Lasky, L. A., Sander, C., Boone, C., Bader, G. D., and Sidhu, S. S. (2008) A specificity map for the PDZ domain family. *PLoS Biol.* 6, e239.
- [11] Tyson R. Shepherd, al. Distinct Ligand Specificity of the Tiam1 and Tiam2 PDZ Domains. *Biochemistry* 2011, 50, 1296-1308. DOI: 10.1021/bi1013613.
- [12] T., Simonson, T., Gaillard, D., Mignon, M., Schmidt am Busch, A., Lopes, N., Amara, S., Polydorides, A., Sedano, K., Druart, G., Archontis. Computational Protein Design: The Proteus Software and Selected Applications. *Journal of Computational Chemistry* 2013, 34, 2472–2484
- [13] Dardel, F., Képès, F.. Bioinformatique Génomique et post-génomique. ISBN: 978-2-7302-0927-4
- [14] Computational design of MHC class II binding epitopes

PDB code	Ligand
1IHJ	TEFCA
1KWA	SYREF
1MFG	LDVPV
1MFL	LDVPV
1N7F	RTYSC
1OBY	NEFYA
1V1T	NEYKV
1W9E	NEFYF
1W9O	NEYVV
1W9Q	NEFAF
1YBO	NEFYA
2EJY	KEYCI
2JOA	RIWWV
2PKU	ESVKI
2PZD	TMFWV
3GCN	--YQF
3GCO	NVYQF
3GDS	NVYYF
3GDU	--YRF
3GDV	--YQF
3HPK	ESVKI
3HPM	ESVKI
3KZE	KEYYA
3OTP	GILQI
3R0H	GQYWV
4GVC	EEFYA
4GVD	EEFYA
4Q2N	WFLDI
4RQZ	LVYQF
4UU5	MERLI

table1. The ligand of class II PDZ (30 in total) collecting from PDB. The first column represents the PDB code and the second column represents the ligand sequences at the last 5 positions.

Group1	Ala, Cys, Ser, Thr	Small
Group2	Glu, Asp, Asn, Gln	Acidic polar
Group3	Hip, Hid, Hie	Basic aromatic
Group4	Ile, Val, Leu, Met	Hydrophobic
Group5	Gly, Pro	No mutable in our system
Group6	Arg, Lys	Basic
Group7	Trp, Phe, Tyr	Aromatic

table2. Amino acid groups and their properties. We classified the residues into 7 groups according to the matrix BLOSUM62.



fexp	Buried residues			Exposed residues			All residues		
	Tiam1	Cask	Mean	Tiam1	Cask	Mean	Tiam1	Cask	Mean
ALA	0,071	0,047	0,059	0,073	0,019	0,046	0,065	0,032	0,048
CYS	0,000	0,030	0,015	0,023	0,001	0,012	0,016	0,013	0,015
SER	0,049	0,151	0,044	0,166	0,158	0,320	0,053	0,236	0,146
THR	0,031	0,044	0,038	0,073	0,079	0,076	0,050	0,056	0,053
GLU	0,060	0,063	0,061	0,114	0,097	0,105	0,093	0,073	0,083
ASP	0,031	0,121	0,040	0,124	0,035	0,084	0,040	0,062	0,046
ASN	0,030	0,007	0,019	0,123	0,309	0,060	0,087	0,328	0,222
GLN	0,001	0,014	0,008	0,051	0,123	0,087	0,032	0,065	0,048
HIS	0,002	0,002	0,013	0,013	0,007	0,007	0,014	0,080	0,008
ILE	0,117	0,197	0,157	0,048	0,035	0,041	0,067	0,098	0,082
VAL	0,132	0,520	0,139	0,572	0,135	0,038	0,072	0,055	0,092
LEU	0,256	0,151	0,204	0,546	0,056	0,144	0,040	0,175	0,319
MET	0,016	0,084	0,050	0,001	0,028	0,014	0,006	0,048	0,027
ARG	0,029	0,088	0,006	0,078	0,018	0,060	0,130	0,222	0,065
LYS	0,059	0,072	0,065	0,083	0,111	0,092	0,101	0,196	0,139
TRP	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000
PHE	0,060	0,118	0,040	0,040	0,050	0,079	0,002	0,013	0,020
TYR	0,058	0,001	0,029	0,015	0,015	0,002	0,009	0,029	0,001
GLY	0,000	0,001	0,000	0,007	0,000	0,004	0,019	0,025	0,018
PRO	0,001	0,001	0,006	0,007	0,003	0,006	0,003	0,004	0,014

table4. Experimental amino acid composition table. The experimental sequences contain 50 Tiam1 homology sequences and 126 hCask homology sequences with identify percentage between 60% and 85%. We separated the buried and exposed residue positions using a buried factor threshold 0.2. In Tiam1, there are 33 buried positions and 51 exposed positions. In Cask, there are 34 buried positions and 40 exposed positions.

	eps8		eps4	
	buried	exposed	buried	exposed
ALA	0,037	0,057	0,05	0,06
CYS	0,028	0	0,027	0
SER	0,056	0,164	0,054	0,171
THR	0,044	0,06	0,04	0,057
GLU	0,066	0,165	0,066	0,167
ASP	0,041	0,137	0,041	0,071
ASN	0,015	0,063	0,017	0,372
GLN	0,015	0,364	0,015	0,057
HIP	0,004	0,043	0,003	0,036
HID	0,002	0,008	0,002	0,007
HIE	0,001	0,058	0,002	0,008
ILE	0,159	0,028	0,156	0,027
VAL	0,099	0,038	0,095	0,034
LEU	0,155	0,515	0,126	0,16
MET	0,102	0,041	0,511	0,042
ARG	0,054	0,186	0,048	0,177
LYS	0,049	0,104	0,241	0,097
PHE	0,05	0,055	0,049	0,06
TYR	0,022	0,072	0,011	0,013
GLY	0	0	0	0
PRO	0	0	0	0

table5. Calculated amino acid composition table. The calculated sequences are from the last iteration for optimization of empirical reference energies.

	MC energies (simulation1)						MC energies (simulation2)					
	Position -2			Position 0			Position -2			Position 0		
	Averaged	Min	Max	Averaged	Min	Max	Averaged	Min	Max	Averaged	Min	Max
A	-26,482	-32,996	-24,157	-25,669	-35,508	-23,547	-34,646	-35,270	-32,862	-28,925	-55,374	-27,142
C	-24,832	-30,825	-22,544	-	-	-	-31,507	-34,788	-30,877	-	-	-
D	-27,387	-32,548	-24,503	-	-	-	-	-	-	-	-	-
E	-23,999	-29,704	-21,768	-	-	-	-31,762	-32,042	-31,678	-	-	-
F	-21,895	-32,214	-19,752	-	-	-	-28,920	-34,675	-27,142	-	-	-
H	-22,924	-31,404	-20,673	-	-	-	-30,352	-55,270	-28,736	-	-	-
K	-23,821	-43,049	-21,992	-	-	-	-31,722	-32,065	-31,532	-	-	-
L	-27,853	-28,920	-27,300	-	-	-	-31,769	-32,137	-31,448	-	-	-
M	-23,247	-32,900	-21,126	-	-	-	-29,633	-33,569	-27,972	-	-	-
N	-24,488	-32,839	-22,028	-	-	-	-31,784	-35,709	-31,151	-	-	-
Q	-24,625	-30,879	-22,344	-	-	-	-31,999	-32,124	-31,833	-	-	-
R	-21,700	-32,283	-19,592	-	-	-	-30,386	-33,421	-28,968	-	-	-
S	-25,666	-31,248	-23,017	-28,473	-38,976	-25,898	-32,149	-37,082	-31,705	-30,244	-32,260	-28,348
T	-25,409	-29,890	-23,509	-58,745	-61,124	-56,354	-34,706	-34,706	-34,706	-	-	-
h	-22,067	-32,717	-19,805	-	-	-	-30,081	-33,883	-28,379	-	-	-
j	-22,313	-32,275	-19,989	-	-	-	-30,300	-35,833	-28,670	-	-	-

table6. proteus energies (folding free energies) from MC by single point mutant explorations of peptide sequences (at position -2 and 0) with Tiam1:Syndecan1 complex. In proteus, the sequences are generated by the highest energies. The system is built using the minimized structure from Xray (4GVD\_A\_D). The complete sequence of peptide is "TKQEEFYA". The MC simulations used frozen protein and eps8. In simulation1, we used the reference energies obtained from extended peptide model. In simulation2, we used the reference energies obtained from empirical model with multiple proteins (Tiam1+Cask).

	MC energies (simulation1)						MC energies (simulation2)					
	Position 0			Position -2			Position 0			Position -2		
	Averaged	Min	Max	Averaged	Min	Max	Averaged	Min	Max	Averaged	Min	Max
A	-407,4	-408,5	-403,1	-362,7	-363,7	-358,8	-434,3	-435,4	-429,9	-436,3	-437,3	-432,4
C	-406,8	-407,6	-403,7	-360,7	-361,5	-356,7	-433,2	-434,1	-430,1	-434,5	-435,3	-430,4
D	-416,9	-417,7	-413,6	-368,1	-369,9	-363,6	-446,0	-446,9	-442,7	-443,2	-445,0	-438,8
E	-419,6	-420,6	-415,9	-362,0	-362,9	-357,6	-447,6	-448,6	-443,8	-437,3	-438,2	-432,9
F	-452,0	-452,9	-448,1	-361,5	-362,6	-357,1	-478,9	-479,8	-474,9	-434,3	-435,4	-429,9
H	-416,7	-417,6	-413,6	-364,2	-365,1	-360,5	-445,0	-445,8	-441,9	-437,6	-438,5	-433,9
h	-413,6	-414,5	-410,2	-363,2	-364,1	-359,1	-441,9	-442,7	-438,5	-437,2	-438,0	-433,1
j	-419,7	-420,5	-416,9	-364,0	-364,8	-361,1	-448,3	-449,1	-445,4	-438,1	-438,9	-435,2
I	-415,8	-416,7	-412,5	-368,8	-369,7	-365,3	-440,8	-441,6	-437,5	-441,4	-442,3	-437,9
K	-423,1	-424,0	-420,0	-361,8	-362,6	-358,6	-450,3	-451,2	-447,1	-436,7	-437,6	-433,6
L	-423,6	-424,5	-420,5	-364,7	-365,6	-361,6	-449,1	-450,0	-446,0	-436,9	-437,7	-433,8
M	-416,4	-417,1	-412,5	-358,8	-359,8	-354,5	-441,6	-442,3	-437,7	-431,1	-432,1	-426,7
N	-414,8	-417,2	-409,8	-363,1	-364,0	-359,3	-443,7	-446,1	-438,6	-437,5	-438,4	-433,7
Q	-419,0	-419,9	-415,3	-362,7	-363,6	-359,0	-446,9	-447,8	-443,3	-437,6	-438,5	-433,8
R	-426,8	-427,7	-422,5	-360,7	-361,7	-356,9	-454,2	-455,1	-449,9	-435,5	-436,5	-431,6
S	-409,4	-410,3	-406,0	-364,0	-364,8	-360,3	-435,1	-436,0	-431,7	-438,1	-438,9	-434,4
T	-413,1	-414,0	-409,3	-363,3	-364,2	-359,4	-438,1	-439,0	-434,3	-437,5	-438,4	-433,6
V	-415,6	-416,4	-412,9	-367,8	-368,6	-363,9	-439,9	-440,8	-437,2	-440,5	-441,3	-436,6
W	-458,9	-459,8	-455,1	-361,9	-362,7	-359,0	-486,4	-487,3	-482,6	-435,1	-435,8	-432,2
Y	-486,2	-487,0	-483,3	-361,1	-362,0	-357,1	-512,9	-513,7	-510,0	-434,1	-435,1	-430,1

table7. proteus energies (folding free energies) from MC by single point mutant explorations of peptide sequences (at position 0 and -2) with Tiam1:Syndecan1 complex. In proteus, the sequences are generated by the highest energies. The system is built using the minimized structure from Xray (4GVD\_A\_D). The complete sequence of peptide is "TKQEEFYA". The MC simulations used flexible protein and eps8. In simulation1, we used the reference energies obtained from extended peptide model. In simulation2, we used the reference energies obtained from empirical model with multiple proteins (Tiam1+Cask).



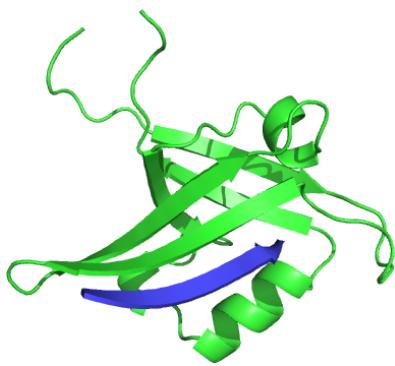


figure1. 3D structure of Tiam1PDZ (in green) binding Syndecan1 peptide (in blue).

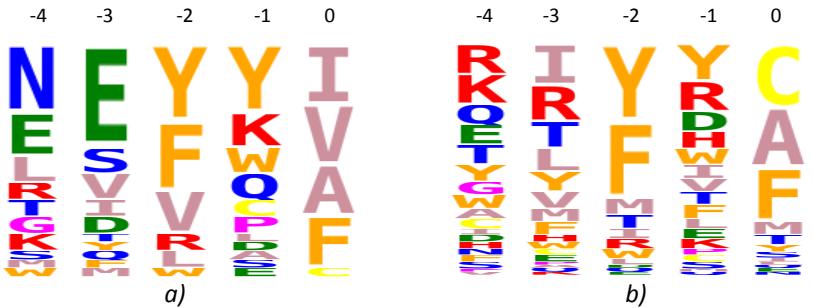


figure2. Sequence logos of the peptide ligands from experiments. a) Ligands of class II PDZ (26 in total) collecting from PDB website. b) Ligands are selected from a Tiam1-binding peptide library (69 in total)<sup>[11]</sup>. The complete sequences of these peptides are “YAA<sub>X<sub>4</sub></sub>X<sub>3</sub>X<sub>2</sub>X<sub>-1</sub>X<sub>0</sub>” , where X represents a given residue type.

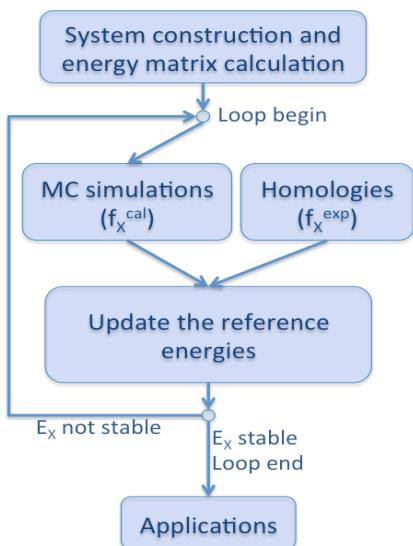


figure3. Flowchart of the general method for the reference energy optimization using empirical model.

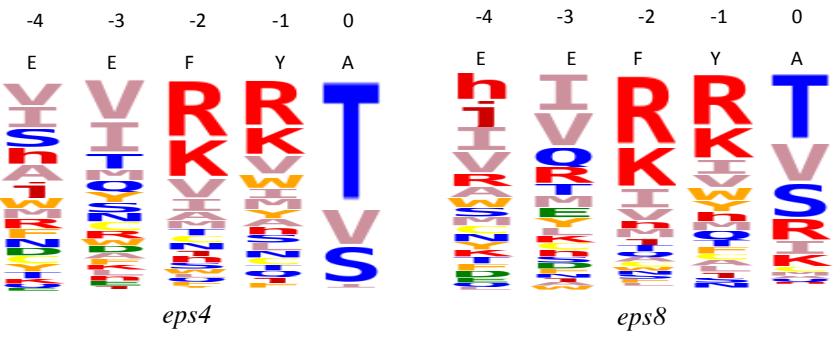
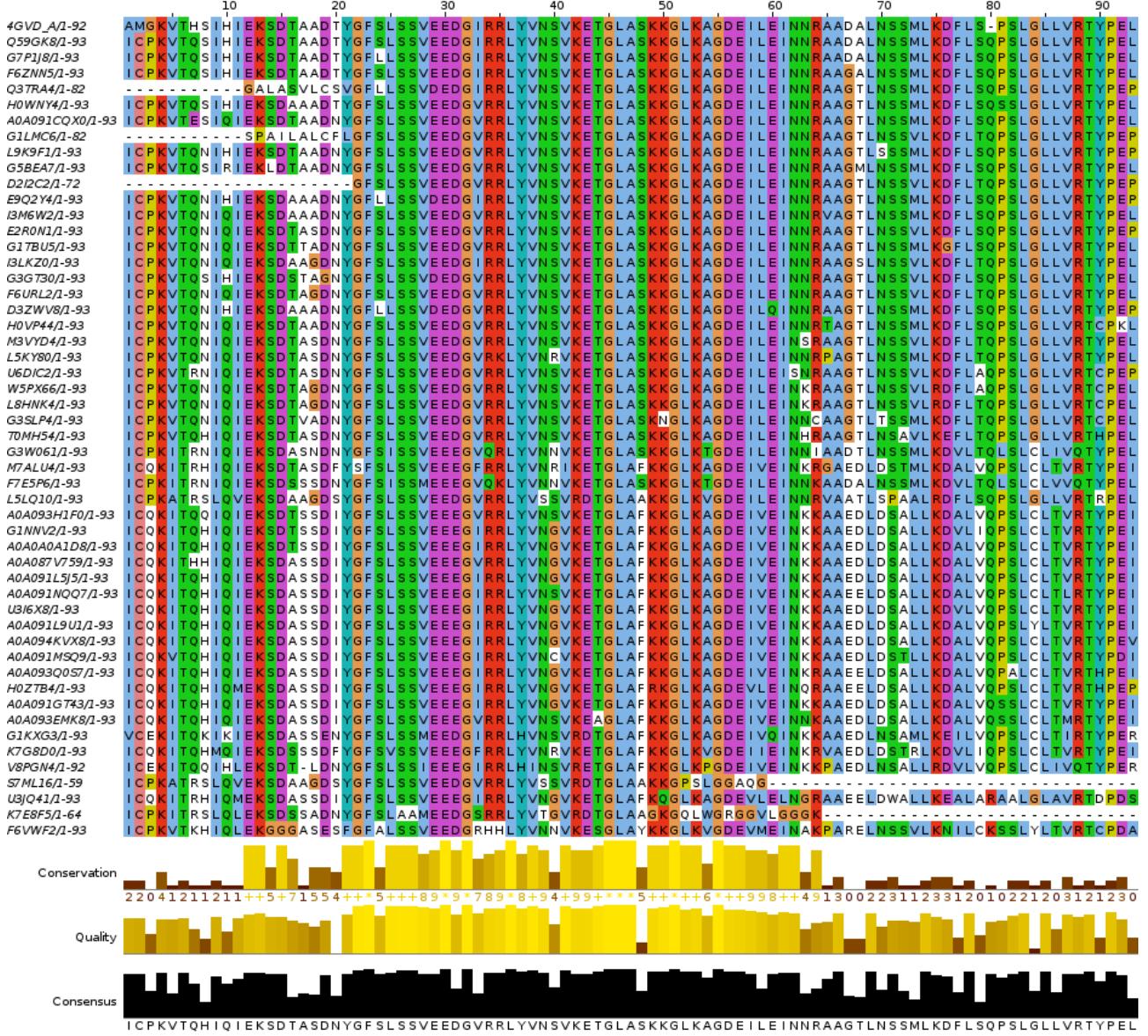


figure4. Sequence logos from MC by single point mutations at the last 5 active positions of peptide using extended peptide model. The complete sequence of peptide is “TKQEEFYA”. The system is built with the extended peptide structure, which is extracted from Xray complex (4GVD\_A\_D) and minimized slightly.



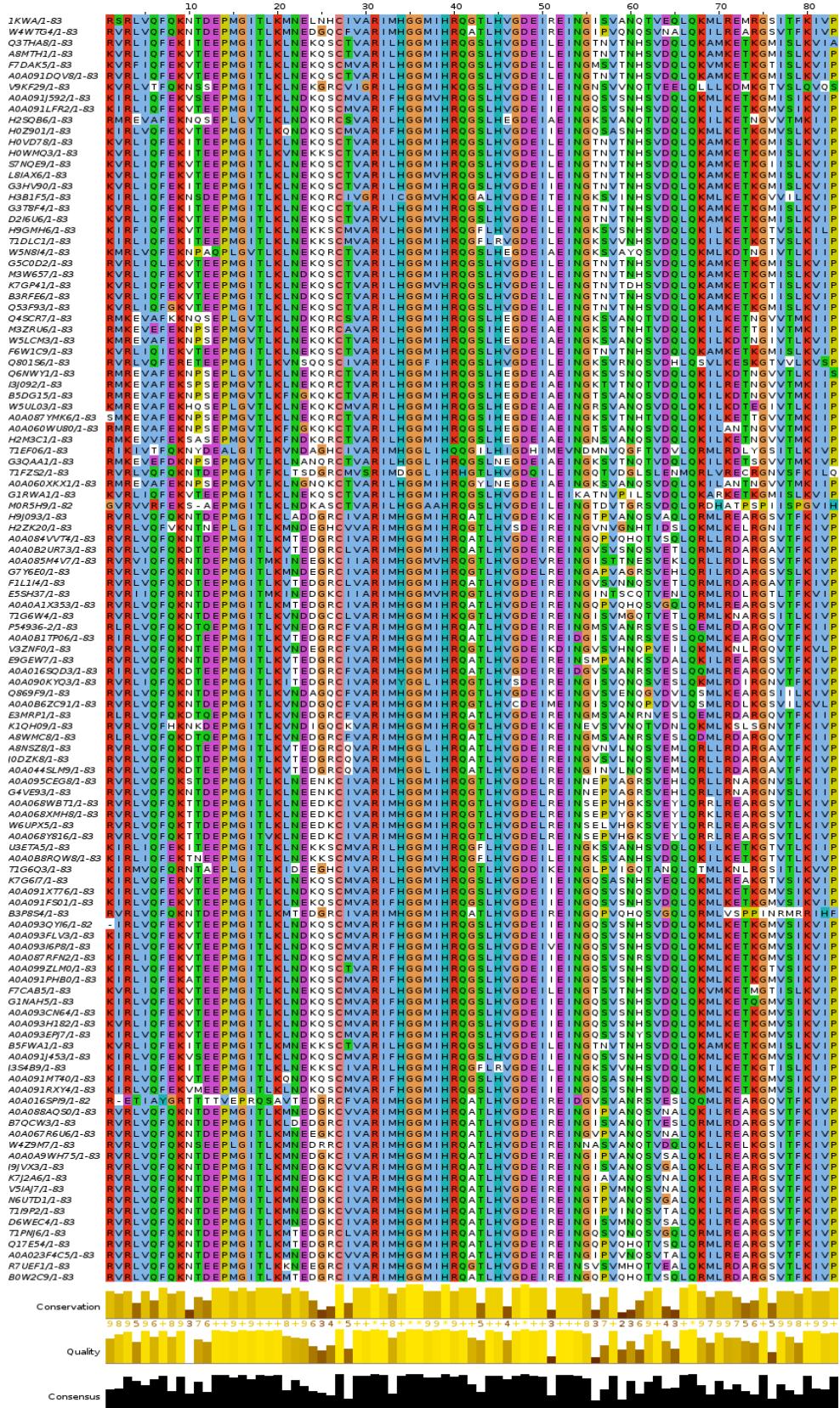


figure6. Sequence alignment with 126 homologues of PDZ hCask protein (PDB: 1KWA\_A). The selected homologues have a percentage of identification between 60% and 85%. The sequences are ordered by the percentage of identification, and the first sequence is the wild type sequence.

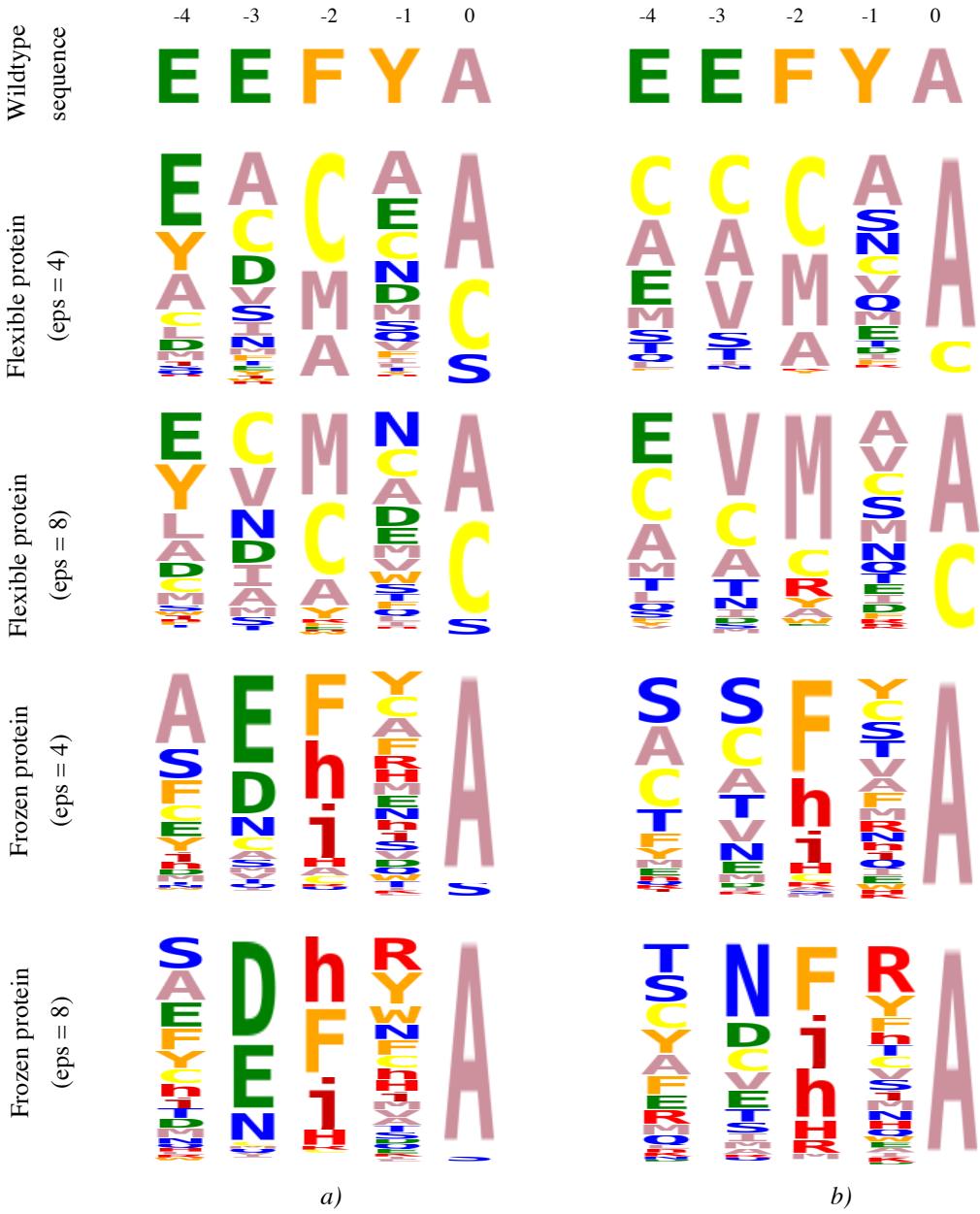


figure7. Sequence logo from MC by mutation explorations of peptide sequences with Tiam1:Syndecan1 complex. The system is built using the minimized structure from Xray (4GVD\_A\_D). The complete sequence of peptide is “TKQEEFYA”. The reference energies are used the energies that obtained from the simulations of the syndecan peptide. For the MC simulations, we considered the protein as flexible or frozen, and used the dielectric constant equals 4 or 8. We explored mutations at the positions -4 to 0 of the peptide ligand. a): The multiple mutation results by 10 million step MC. b): The single point mutation results by 5 million step MC at each position.

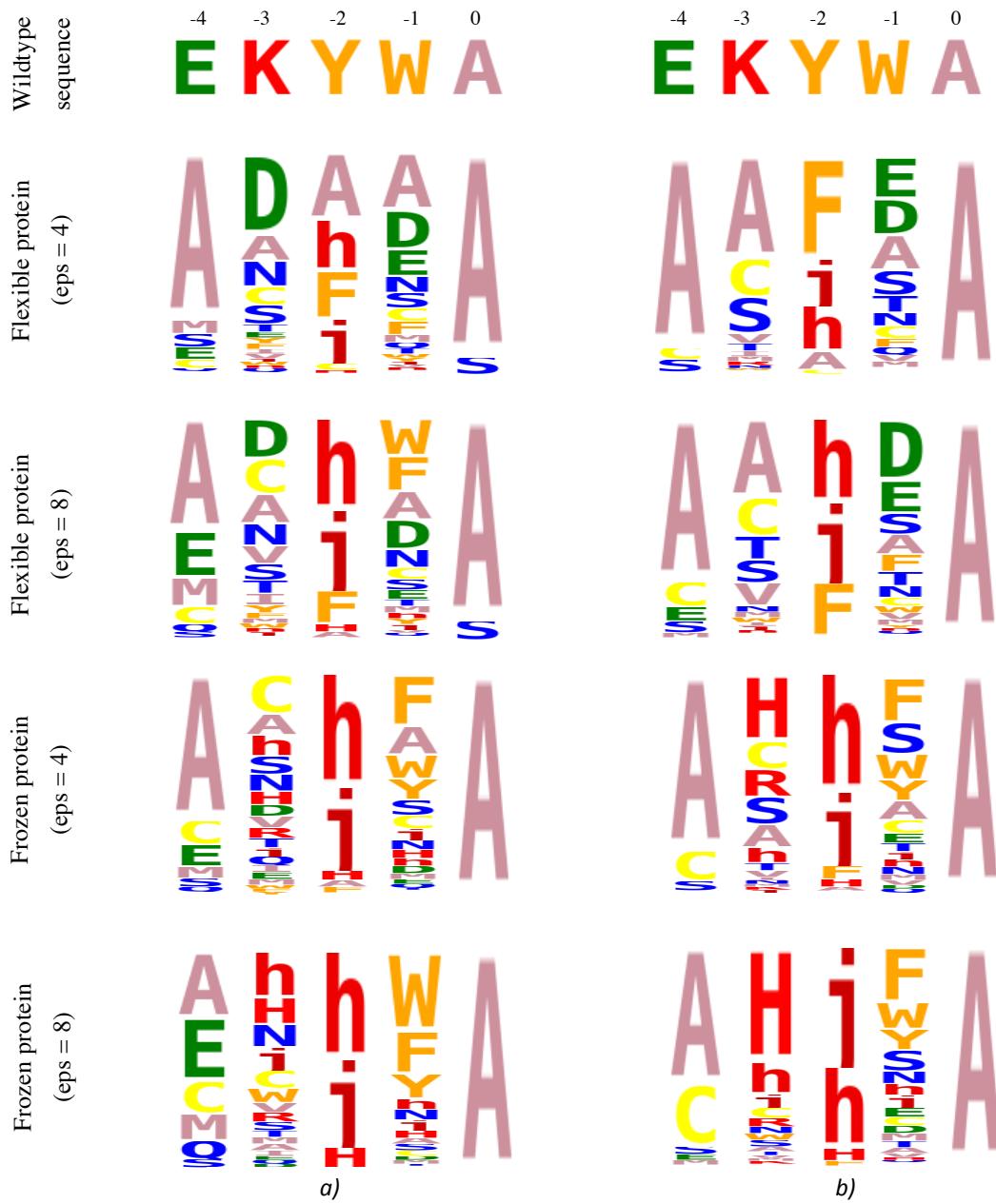


figure8. Sequence logo from MC simulation of Tiam1:peptide. The system is built using the structure obtained from 20ns MD. “Wildtype sequence”(above) is from the Tiam1-binding peptide library. The complete sequence is “YAAEKYWA”<sup>[11]</sup>. The reference energies are used the energies that obtained from the simulations of the syndecan peptide. For the MC simulations, we considered the protein as flexible or frozen, and used the dielectric constant equals 4 or 8. We explored mutations at the positions -4 to 0 of the peptide ligand. a): The multiple mutation results by 10 million step MC. b): The single point mutation results by 5 million step MC at each position.

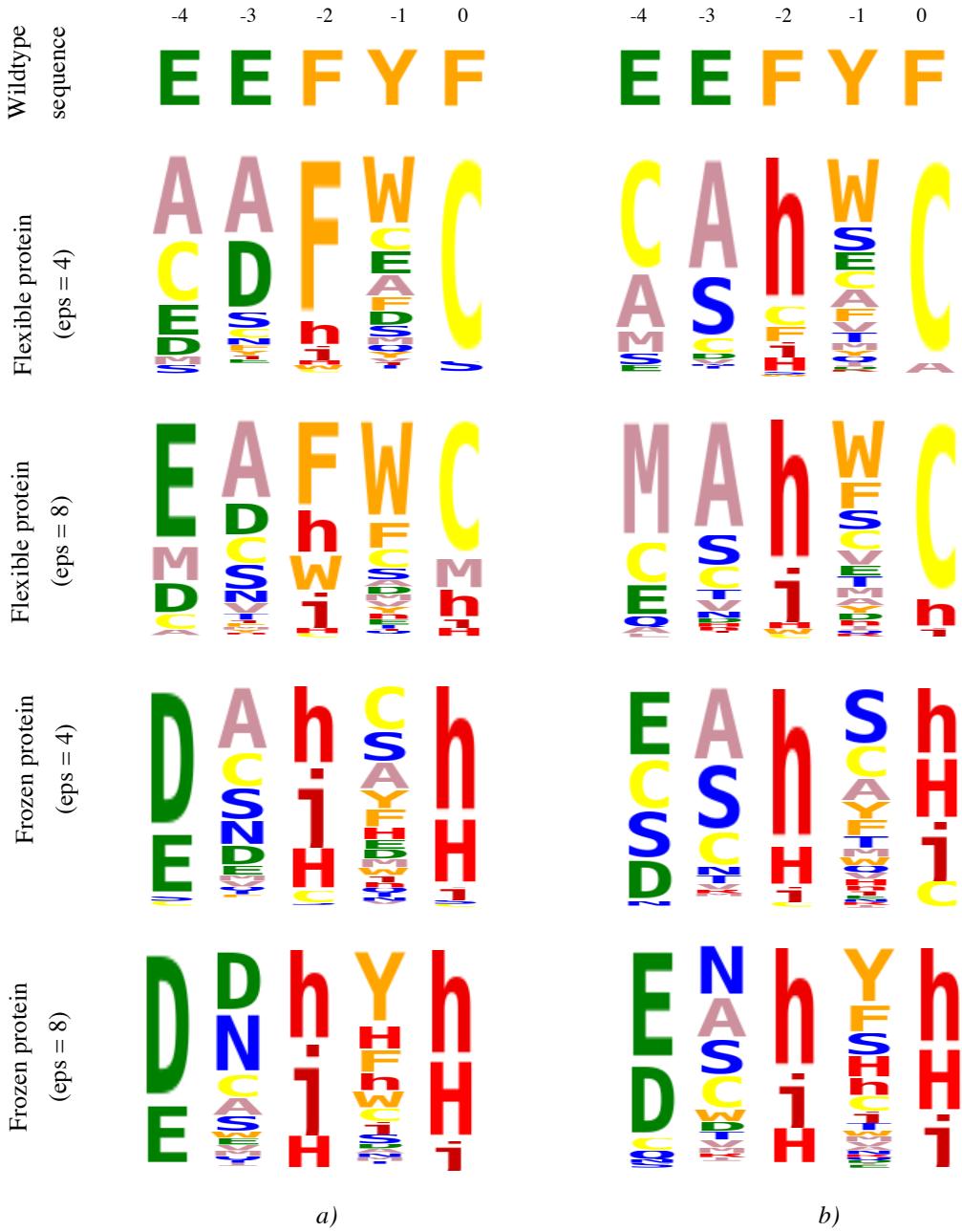
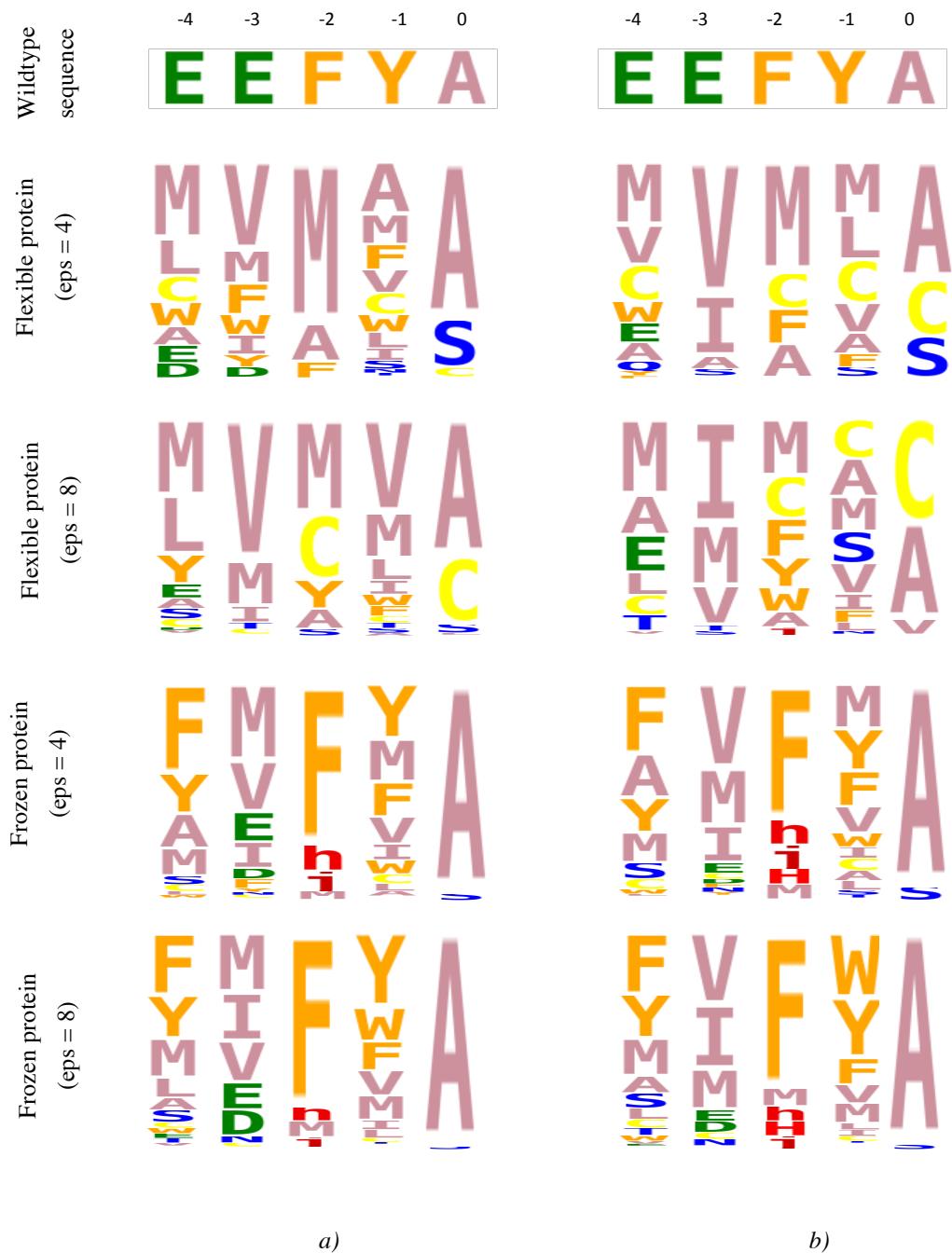


figure9. Sequence logo from MC simulation of Tiam1:peptide. The system is built using the averaged structure from 100ns MD. The “Wildtype sequence” (above) is from a mutant of Syndecan1 peptide. The complete sequence is “TKQEEFYF”. The reference energies are used the energies that obtained from the simulations of the syndecan peptide. For the MC simulations, we considered the protein as flexible or frozen, and the dielectric constant equals to 8 or 4. We explored mutations at the positions -4 to 0 of the peptide ligand. a): The multiple mutation results by 10 million step MC. b): The single point mutation results by 5 million step MC at each position.



*figure10.* Sequence logo from MC simulation of Tiam1:peptide. “wildtype” is the Syndecan1 peptide from the complex of the PDZ domain of Tiam1; the complete sequence is “TKQEEFYA”. The reference energies are used the reference energies, which is optimized using multiple proteins (Cask and Tiam1) with 1/5 active positions. For the MC simulations, we considered the protein as flexible or frozen, and used the dielectric constant equals 4 or 8. We explored mutations at the positions -4 to 0 of the peptide ligand. *a)*: The multiple mutation results by 10 million step MC. *b)*: The single point mutation results by 5 million step MC at each position.

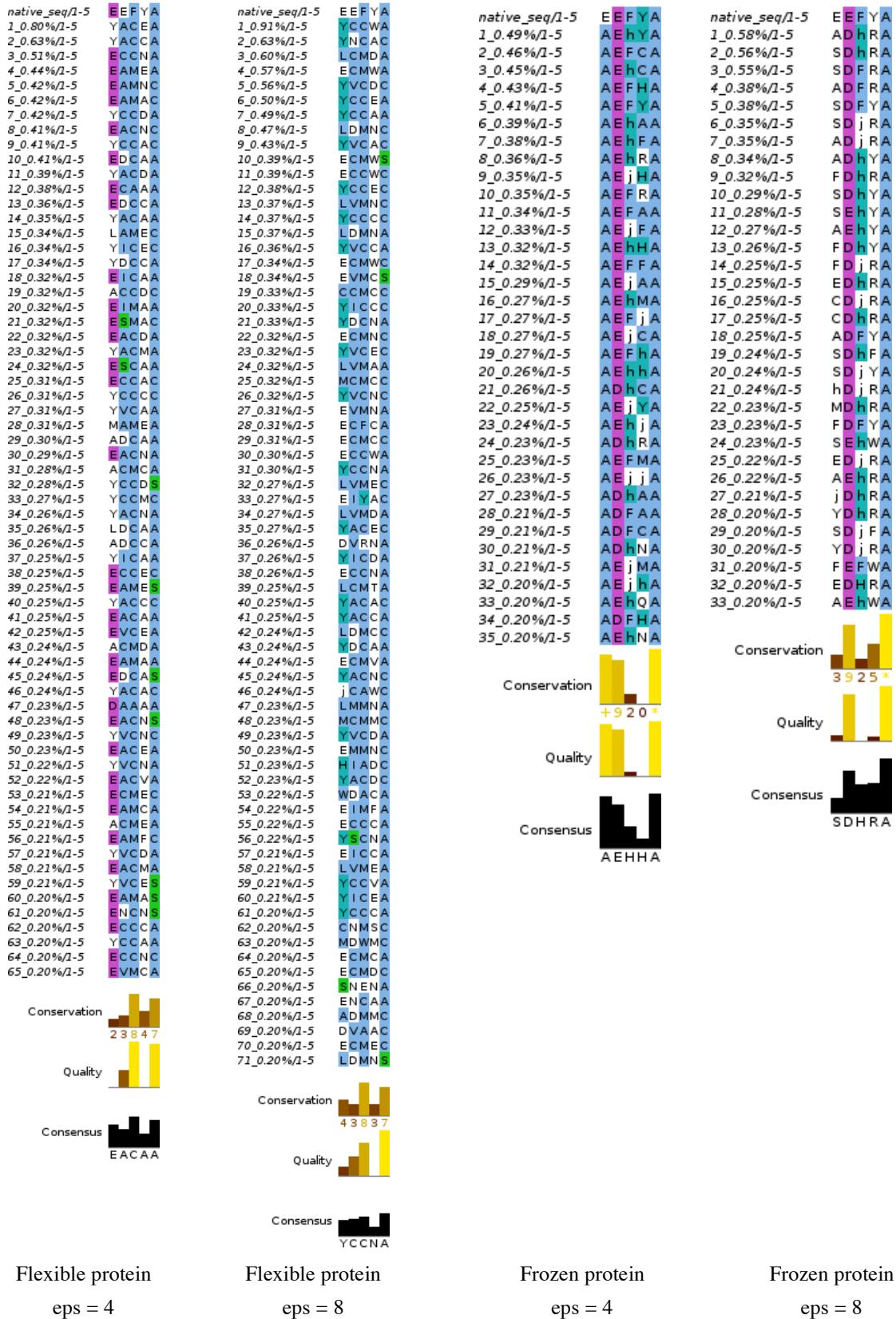


figure11. Sequence alignments from MC by multiple mutation explorations of peptide sequences with Tiam1:Syndecan1 complex. The system is built using the minimized structure from Xray (4GVD\_A\_D). The complete sequence of peptide is “TKQEEFYA”. For the MC simulations, we considered the protein as flexible or frozen, and used the dielectric constant equals 4 or 8. The reference energies are used the delta mutant energies from extended peptide. We explored mutations at the positions 0 to -4 of the peptide ligand. The occurrence threshold for selecting sequences is 0.2%. The first lines are the native sequence and the title of the other lines contains “seqid\_occurrences/1-5”. The sequences are ordered by their occurrences.

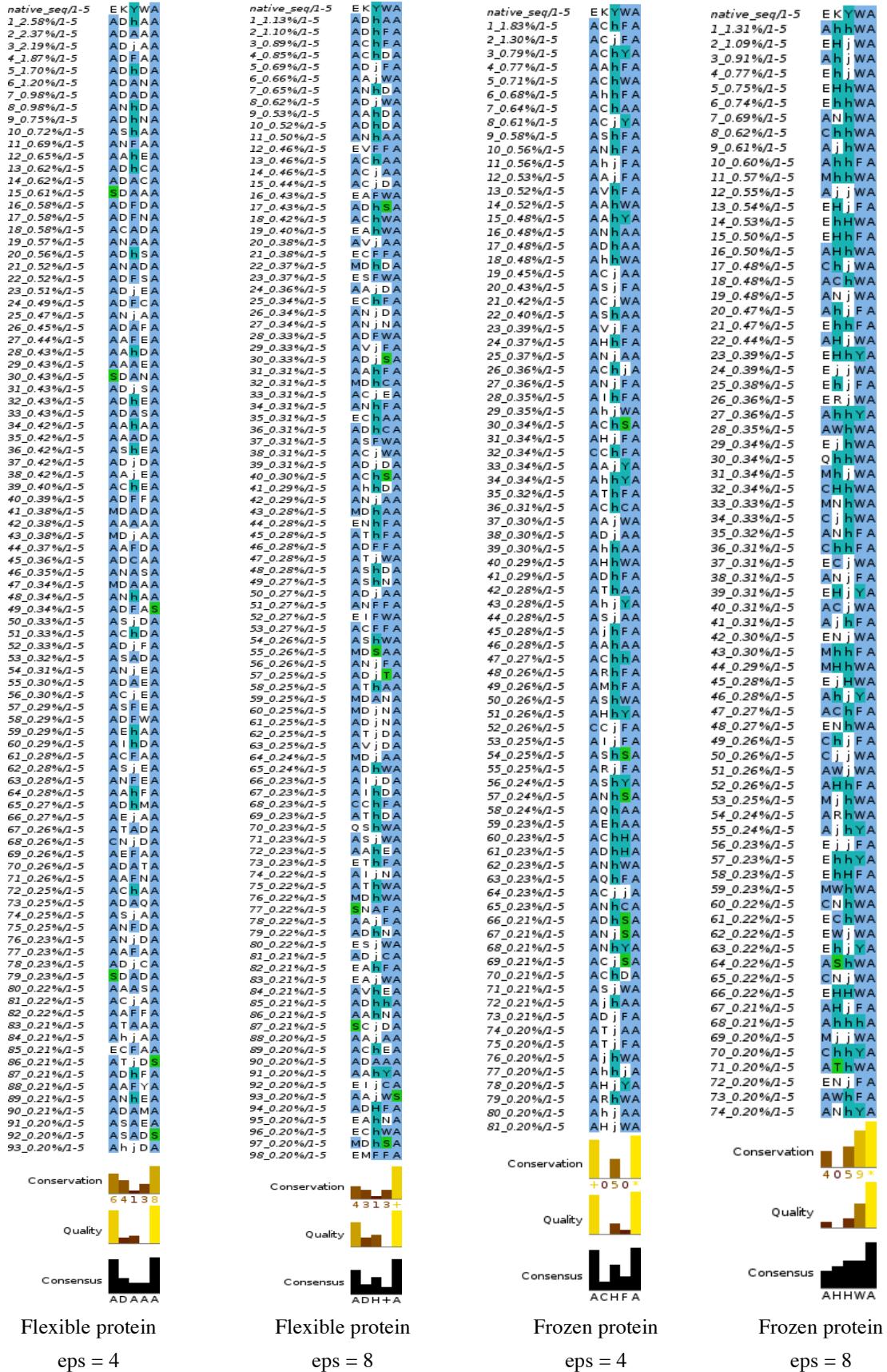


figure12. Sequence alignments from MC by multiple mutation explorations of peptide sequences with Tiam1:peptide complex. The system is built using the structure obtained from 20ns MD. The peptide is selected from the Tiam1-binding peptide library with the complete sequence "YAAEKYWA"<sup>[11]</sup>. For the MC simulations, we considered the protein as flexible or frozen, and used the dielectric constant equals 4 or 8. The reference energies are used the delta mutant energies from extended peptide. We explored mutations at the positions 0 to -4 of the peptide ligand. The occurrence threshold to select sequences is 0.2%. The first lines are the native sequence and the title of the other lines contains "seqid\_occurrences/1-5". The sequences are ordered by their occurrences.

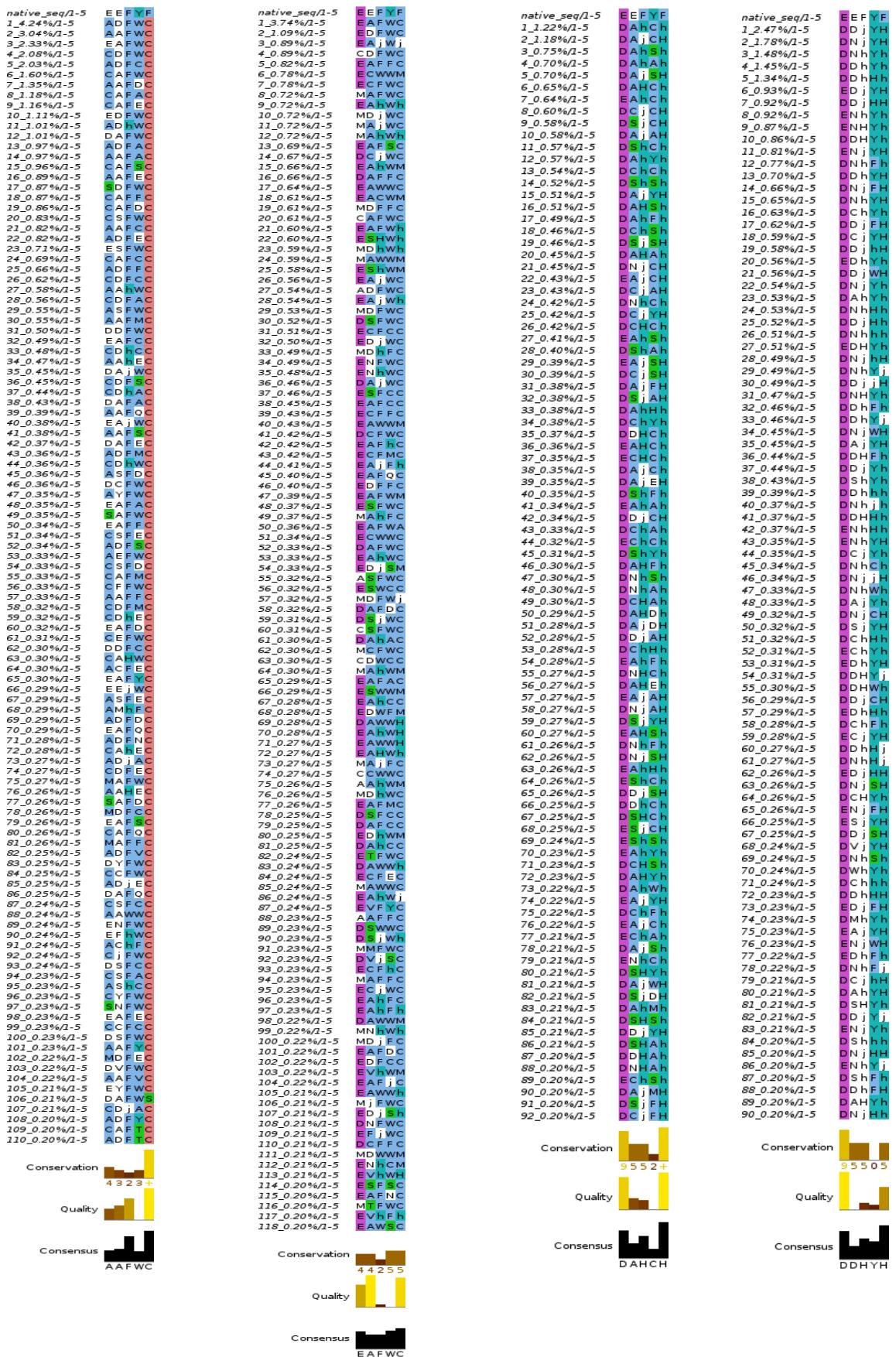


figure13. Sequence alignments from MC by multiple mutation explorations of peptide sequences with Tiam1:peptide complex. The system is built using the averaged structure from 100ns MD. The peptide is from a mutant of Syndecan1 peptide with the complete sequence “TKQEEFYF”. For the MC simulations, we considered the protein as flexible or frozen, and used the dielectric constant equals 4 or 8. The reference energies are used the delta mutant energies from extended peptide. We explored mutations at the positions 0 to -4 of the peptide ligand. The occurrence threshold for selecting sequences is 0.2%. The first lines are the native sequence and the title of the other lines contains “seqid\_occurrences/1-5”. The sequences are ordered by their occurrences.

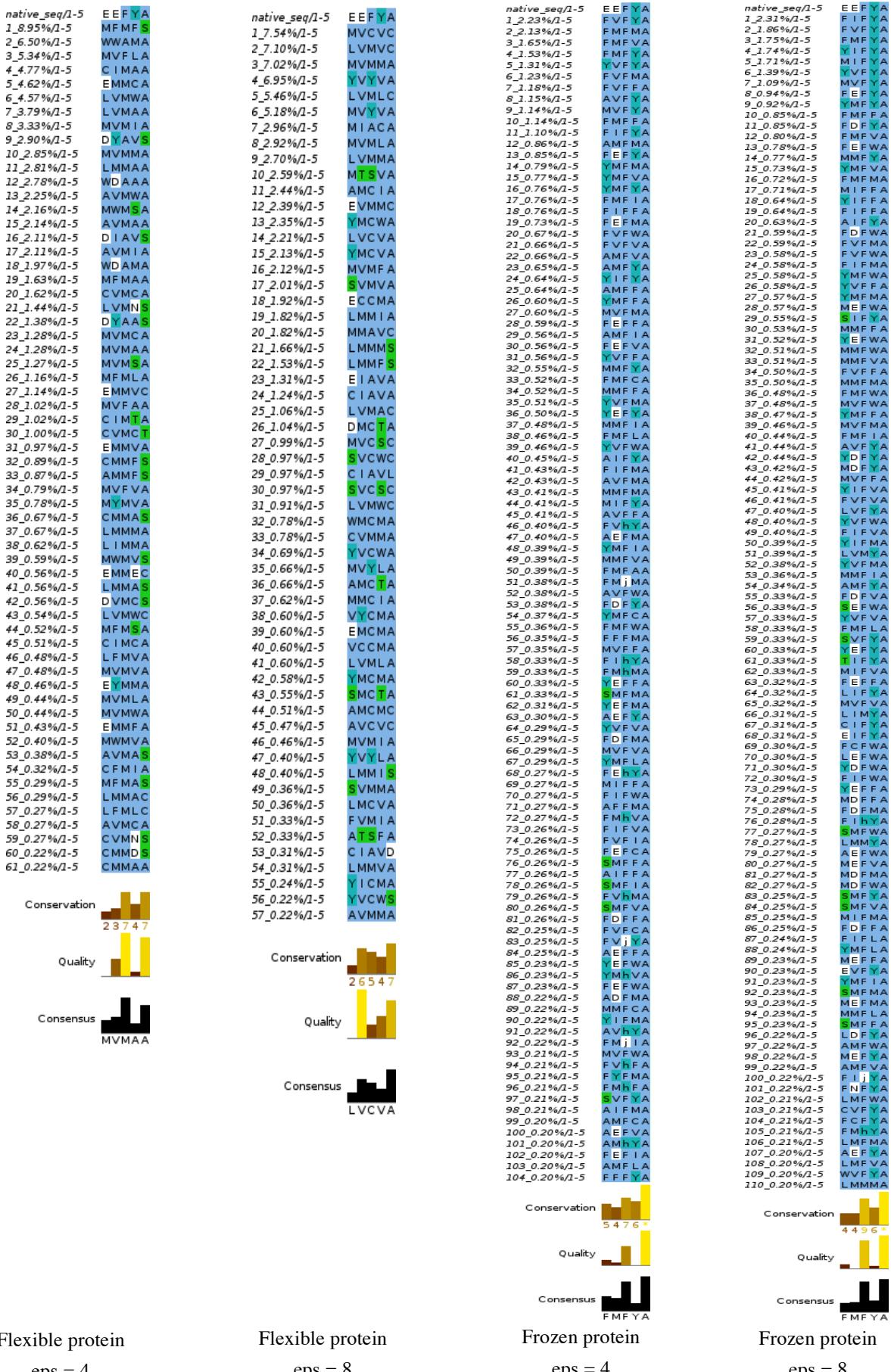
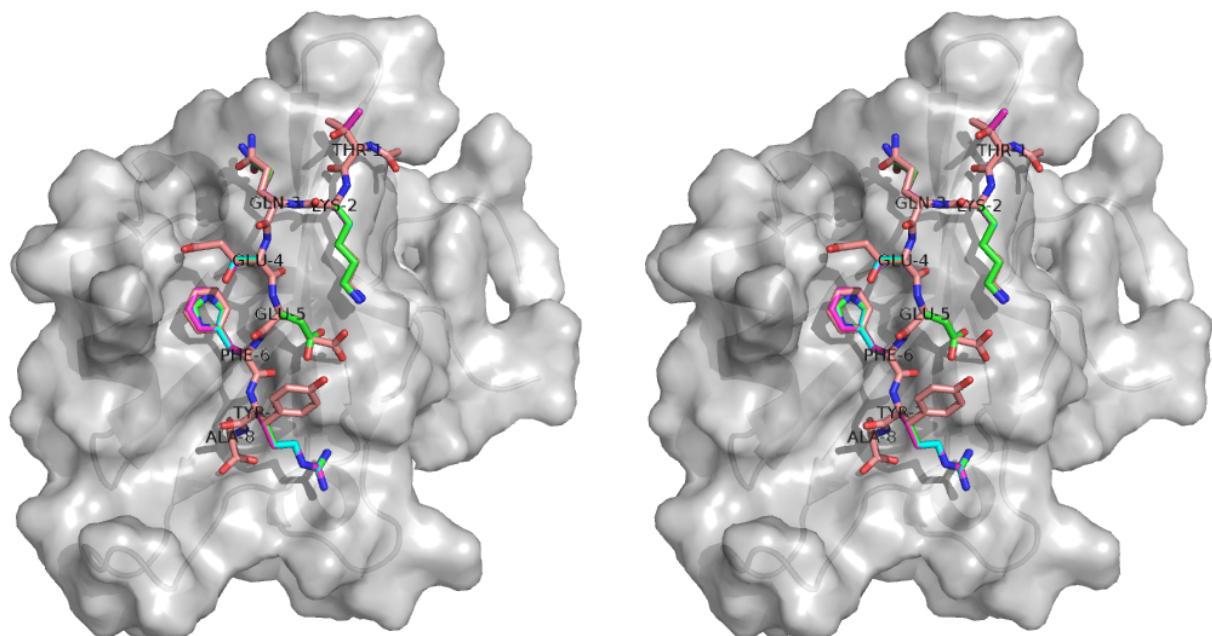
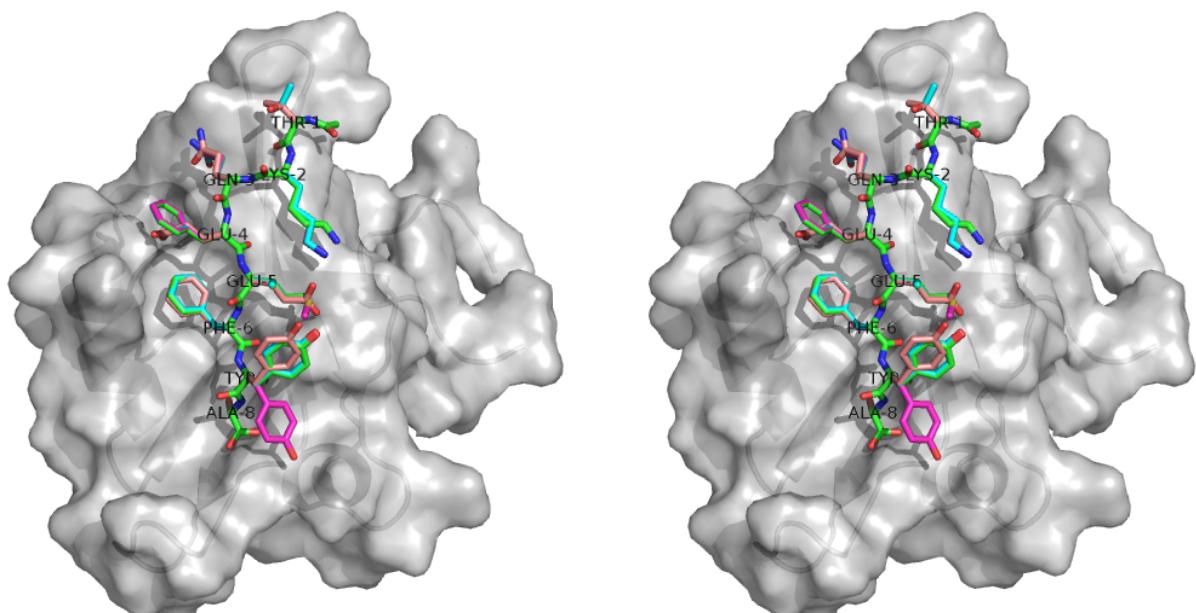


figure14. Sequence alignments from MC by multiple mutation explorations of peptide sequences with Tiam1:Syndecan1 complex. The system is built using the minimized structure from Xray (4GVD\_A\_D). The complete sequence of peptide is “TKQEEFYA”. For the MC simulations, we considered the protein as flexible or frozen, and used  $\epsilon = 4$  or 8. The reference energies are optimized using multiple proteins (Tiam1+Cask) with 1/5 active positions. We explored mutations at the positions 0 to -4 of the peptide ligand. The occurrence threshold for selecting sequences is 0.2%. The first lines are the native sequence and the title of the other lines contains “seqid\_occurrences/1-5”. The sequences are ordered by their occurrences.



a)



b)

*figure15. Cross-eye stereo of the reconstructed structures and of the Xray structure (slightly minimized). The sequences of reconstructed structures are obtained from MC with Tiam1:Syndecan1 complex by multiple mutations, using frozen protein and eps8, which have top3 populations. The protein is represented in grey. The experimental structure of peptide is represented in salmon with the label. The reconstructed structures are represented in bleu, green and magenta. a) The MC simulation is used the reference energies from extended peptide model. The obtained sequences with top3 populations are: “TKQADhRA”, “TKQSDhRA” and “TKQSDFRA”. b) The MC simulation is used the reference energies from empirical model – version2 (using multiple protein Cask+Tiam1). The obtained sequences with top3 populations are: “TKQFIFYA”, “TKQFVFY” and “TKQFMFYA”.*

## APPENDIX

The directories in the path `/home/xingyu/` are for the simulations with different structures:

`/sdc_Proteus/` : Syndecan peptide alone (with first 3 inactive positions and last 5 active positions).

`/PDZ_complexe_TKQ/` and `/PDZ_complexe_TKQpf/` : Tiam1:Syndecan complex with flexible or frozen (pf) protein and the last 5 position flexible in peptide.

`/PDZ_complexe_YAA/` and `/PDZ_complexe_YAApf/` : Tiam1:peptide complex with flexible or frozen (pf) protein and the last 5 position flexible in peptide. The peptide is from the Tiam1-binding peptide library with sequence “YAAEKYWA”<sup>[11]</sup>.

`/PDZc_A0F/` and `/PDZc_A0F_pf/` : Tiam1:peptide complex with flexible or frozen (pf) protein and the last 5 position flexible in peptide. The peptide is from a mutant of Syndecan1 peptide with sequence “TKQEEFYF”.

`/PDZ_Protein_act/` and `/Cask_Pro_act/` : Tiam1 or Cask protein alone with all positions as active.

`/Proteus_2.0.1/`, Proteus package with some modified scripts:

1. The topology for N3 are added in the files `.str` and `.inp` :

```
parameters
DIHE N3    CT   C    N    MULT 4 0.0000000 4 0.0 ! four amplitudes and
                                         0.55000000 3 180.0 ! phases for psi
                                         1.58000000 2 180.0
                                         0.45000000 1 180.0
end
```

2. A new version of `analyze_proteus_sequences.pl`

The reference energy optimization with Tiam1 and Cask can be achieved by scripts in cluster with the path (launch the job to 10 nodes in cluster):

```
/work/xingyu/PDZ_Protein_act/proteus/optEner_TC/
/work/xingyu/Cask_Pro_act/proteus/optEner_TC/
```

or by a version (calculate in local machine) with the path:

```
/home/xingyu/PDZ_Protein_act/proteus/optEner_v2/
```

A backup in the directory:

```
/home/xingyu/testcase/optioner_Tiam1_Cask/
```

We have scanned the 20 point mutants at position 0 and -2 of peptide with MC simulation in Tiam1:syndecan (flexible protein and eps8). The directory is:

```
/home/xingyu/PDZ_complexe_TKQ/proteus/run_impo.sh
```

A backup in the directory:

```
/home/xingyu/testcase/pt_mut_scan/
```

*script1*. The shell script for the reference energies optimization (in cluster).

```
#!/bin/bash
export CPD="/bioinfo/CPD/Proteus_2.0.1"
export MATRIX="/work/xingyu/PDZ_Protein_act/matrix"
export PROTEUSEXE="/work/xingyu/proteus/proteus"
export MYPATH="/work/xingyu/PDZ_Protein_act/proteus/optEner_TC"
export CASKPATH="/work/xingyu/Cask_Pro_act/proteus/optEner_TC"
export NODE=node$1
export MYPATHNODE=$MYPATH/$NODE
export i=`grep starting $MYPATH/iteration.log | awk '{print $3}'`
```

```

#copy the reference energies to the current directory
cp $MYPATH/refener_oldb.dat $MYPATHNODE/refener_oldb.dat
cp $MYPATH/refener.olde.dat $MYPATHNODE/refener.olde.dat

export prot=$1
echo iteration $i >> $MYPATHNODE.log
echo "protein: $prot" >> $MYPATHNODE.log
# select the active positions
awk j="$prot'{if(NR%5 ==j%5) printf "%s,", $1}' $MATRIX/position_list.dat >
$MYPATHNODE/pos_act.dat

# Write the option Space_Constraints and Ref_Ener in the configure file of MC
cp $MYPATHNODE/MONTECARLO.conf $MYPATHNODE/MONTECARL01.conf
python $MYPATHNODE/confMC.py $NODE

echo "<Space_Constraints>" >> $MYPATHNODE/MONTECARL01.conf
sort -k1 $MYPATHNODE/space_constraints.dat >> $MYPATHNODE/MONTECARL01.conf
echo "</Space_Constraints>" >> $MYPATHNODE/MONTECARL01.conf

echo "<Ref_Ener>" >> $MYPATHNODE/MONTECARL01.conf
cat $MYPATHNODE/refener_oldb.dat >> $MYPATHNODE/MONTECARL01.conf
sort -k2 $MYPATHNODE/ref_ener.dat >> $MYPATHNODE/MONTECARL01.conf
echo "</Ref_Ener>" >> $MYPATHNODE/MONTECARL01.conf

# Active positions in array
active_resid_array=( `perl -nale 'if($F[1]eq"active"){push@a,$F[0]}END{print
join(q( ),sort{$a<=>$b}@a)}' $MATRIX/position_list.dat` )
# Number of active positions
nbactive=${#active_resid_array[@]}
echo "number of active positions: $nbactive" >> $MYPATHNODE.log
# Make a separated active resid list
active_resid_list=`perl -nale 'if($F[1]eq"active"){push@a,$F[0]}END{print
join(q(),sort{$a<=>$b}@a)}' $MATRIX/position_list.dat`

# Run proteus multiples mutations
$PROTEUSEXE < $MYPATHNODE/MONTECARL01.conf >& $MYPATHNODE/MONTECARL0.log
echo Finished MONTECARL0 of fixed backbone $NODE >> $MYPATHNODE.log
# Select the sequences with temperature 0.592
grep '0.592$' $MYPATHNODE/output.ener_* > $MYPATHNODE/output.ener
cat $MYPATHNODE/proteus_mult.seq_* > $MYPATHNODE/proteus_mult1.seq
python $MYPATHNODE/selcSeq.py $NODE

# Proteus post-processing
$PROTEUSEXE < $MYPATHNODE/POSTPROCESS.conf >& $MYPATHNODE/POSTPROCESS.log
echo Finished POSTPROCESS of fixed backbone $NODE >> $MYPATHNODE.log

```

```

# Analysis of sequences
$MYPATH/analyze_proteus_sequences.pl $MYPATHNODE/proteus_mult.seq
$MYPATHNODE/proteus_mult.rich $active_resid_list 50000000 >
$MYPATHNODE/proteus_mult.dat
cp $MYPATHNODE/proteus_mult.dat $MYPATH/proteus_mult_"$NODE".dat
mv $MYPATHNODE/proteus_mult.seq $MYPATHNODE/proteus_mult_$i.seq
# write the log file
echo $NODE finished > $MYPATH/log/ite_"$NODE"_"$i".log
# check if all 10 separated MC simulation have done
n=`ls $MYPATH/log/*.log | wc -l`
if [ $n -eq 10 ]; then
    for f in $MYPATH/proteus_mult_node1.dat \
        $MYPATH/proteus_mult_node2.dat \
        $MYPATH/proteus_mult_node3.dat \
        $MYPATH/proteus_mult_node4.dat \
        $MYPATH/proteus_mult_node5.dat \
        $CASKPATH/proteus_mult_node1.dat \
        $CASKPATH/proteus_mult_node2.dat \
        $CASKPATH/proteus_mult_node5.dat \
        $CASKPATH/proteus_mult_node4.dat \
        $CASKPATH/proteus_mult_node5.dat ; do
        cat $f >> $MYPATH/proteus_mult.dat
    done
    for f in $MYPATH/log/*.log ; do
        cat $f >> $MYPATH/iteration_$i.log
    done
    rm $MYPATH/log/*.log
# Update refEner with python:
python $MYPATH/refenerOpt1.py $MYPATH/proteus_mult.dat
mv $MYPATH/proteus_mult.dat $MYPATH/proteus_mult_$i.dat
cp $MYPATH/refener_newb.dat $MYPATH/refburied_$i.dat
cp $MYPATH/refener_newe.dat $MYPATH/refexposed_$i.dat
cp $MYPATH/refener_newb.dat $MYPATH/refener_oldb.dat
cp $MYPATH/refener_newe.dat $MYPATH/refener.olde.dat
cp $MYPATH/refener_newb.dat $CASKPATH/refener_oldb.dat
cp $MYPATH/refener_newe.dat $CASKPATH/refener.olde.dat
cp $MYPATH/dat/refdicb.dat $CASKPATH/dat/refdicb.dat
cp $MYPATH/dat/refdice.dat $CASKPATH/dat/refdice.dat
mv $MYPATH/freqCalc.dat $MYPATH/freqCalc_$i.dat
mv $MYPATH/freqCalb.dat $MYPATH/freqCalb_$i.dat

let i=i+1
echo "starting iteration $i" > $MYPATH/iteration.log
ssh master0 $MYPATH/optiEref.sh
fi

```

## **RESUME**

Nous travaillions avec le PDZ domaine de la protéine Tiam1 (*T-cell lymphoma invasion and metastasis gene 1*). Nous voulons tester la méthode Dessin Computationnel de Protéine (DCP) en utilisant différents paramètres : les deux constants diélectriques (4 et 8) et les différents énergies des références ( $E_x$ ). Pour calculer les énergies des références, nous avons développé deux méthodes différents : La première méthode utilise un modèle de peptide déplié. Nous utilisons des Monte Carlo simulations pour calculer les énergies libres des mutants simples par rapport à Ala. La deuxième méthode utilise un modèle empirique, qui va améliorer les  $E_x$  itérativement. La boucle va s'arrêter jusqu'à toutes les compositions des acides aminés calculés convergent vers les expérimentales valeurs. Nous avons réalisé cette méthode en utilisant les protéines multiples (de Tiam1 et de Cask), en regroupant des acides aminés similaires ensemble, et en appliquant la conception de « position – dépendante ». Nous avons utilisé ces  $E_x$  améliorées pour explorer des mutations du peptide syndecan qui est associé avec Tiam1. Les résultats sont satisfaisants en comparant avec des séquences expérimentaux.

## **ABSTRACT**

We work with the PDZ domain of Tiam1. We want to test the performance of Computational Protein Design (CPD) method using different parameters including two dielectric constants (4 and 8) and different reference energies ( $E_x$ ). We used two methods to calculate the reference energies: In the first method, we use an unfolded peptide model and calculate the free energies of all single point mutants relative to Ala by Monte Carlo (MC) simulations. The second method uses an empirical model to iteratively update the  $E_x$ , until the amino acid compositions of calculated sequences are convergent to the experimental values. To implement this empirical model, we use multiple proteins (Tiam1 and Cask), regroup the similar amino acid types and apply position-dependent conception. We applied the improved  $E_x$  to explore mutations of the syndecan peptide bound to Tiam1. We have got satisfied results comparing with the experimental sequences.