



**HAL**  
open science

## Functional Motions Modulating VanA Ligand Binding Unraveled by Self-Organizing Maps

Guillaume Bouvier, Nathalie Duclert-Savatier, Nathan Desdouits, Djalal Meziane-Cherif, Arnaud Blondel, Patrice Courvalin, Michael Nilges, Thérèse Malliavin

► **To cite this version:**

Guillaume Bouvier, Nathalie Duclert-Savatier, Nathan Desdouits, Djalal Meziane-Cherif, Arnaud Blondel, et al.. Functional Motions Modulating VanA Ligand Binding Unraveled by Self-Organizing Maps. *Journal of Chemical Information and Modeling*, 2014, 54 (1), pp.289-301. 10.1021/ci400354b .  
pasteur-02510864

**HAL Id: pasteur-02510864**

**<https://pasteur.hal.science/pasteur-02510864>**

Submitted on 7 Apr 2020

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1 Functional motions modulating VanA ligand binding  
2 unraveled by self-organizing maps

3  
4 November 6, 2013

5 Guillaume Bouvier<sup>1\*</sup>, Nathalie Duclert-Savatier<sup>1\*</sup>, Nathan Desdouits<sup>1</sup>, Djalal Meziane-  
6 Cherif<sup>2</sup>, Arnaud Blondel<sup>1</sup>, Patrice Courvalin<sup>2</sup>, Michael Nilges<sup>1</sup>, Thérèse E. Malliavin<sup>1</sup>

7  
8 (1) Institut Pasteur, Unité de Bioinformatique Structurale; CNRS UMR 3825; Département  
9 de Biologie Structurale et Chimie; 25, rue du Dr Roux, 75015 Paris, France

10 (2) Institut Pasteur, Unité des Agents Antibactériens; 25, rue du Dr Roux, 75015 Paris,  
11 France

12  
13 (\*) contributed equally to the work.

14 Corresponding author : Thérèse E. Malliavin; Unité de Bioinformatique Structurale;  
15 CNRS UMR 3825; Institut Pasteur, 25, rue du Dr Roux, 75015 Paris, France; E-mail:  
16 terez@pasteur.fr; Phone: +33 1 40 61 34 75

17  
18 Keywords : D-alanyl:D-lactate ligase, D-alanyl:D-alanine ligase, antibiotic resistance,  
19 vancomycin, molecular dynamic simulation, docking, classification, self organizing map.

## 21 **1 Abstract**

22 The VanA D-Ala:D-Lac ligase is a key enzyme in the emergence of high level resistance  
23 to vancomycin in *Enterococcus* species and Methicillin-Resistant *Staphylococcus aureus*.  
24 It catalyzes the formation of D-Ala-D-Lac, a surrogate peptidoglycan precursor with low  
25 affinity for vancomycin, that can replace D-Ala-D-Ala, which is subject to sequestration  
26 by vancomycin. Therefore, VanA appears as an attractive target for the design of new  
27 antibacterials to overcome resistance.

28 The catalytic site of VanA is delimited by three domains and closed by an  $\omega$ -loop  
29 upon enzymatic reaction. The aim of the present work was: (i) to investigate the con-  
30 formational transition of VanA associated to the opening of its  $\omega$ -loop; and (ii), to relate  
31 this transition with the substrates or products binding propensities. Molecular dynamics  
32 trajectories of the VanA ligase of *Enterococcus faecium* with or without a disulfide bridge  
33 distant from the catalytic site, revealed differences in the  $\omega$ -loop conformations with a  
34 slight opening. Conformations were clustered with an original machine learning method,  
35 based on self-organizing maps (SOM), which revealed four distinct conformational basins.  
36 Several ligands related to substrates, intermediates or products were docked to SOM rep-  
37 resentative conformations with the DOCK 6.5 program. Classification of ligand docking  
38 poses, also performed with SOMs, clearly distinguished ligand functional classes: sub-  
39 strates, reaction intermediates and product. This result illustrates the acuity of the SOM  
40 classification and supports the quality of the DOCK program poses. The protein-ligand  
41 interaction features for the different classes of poses will guide the search and design of  
42 novel inhibitors.

## 2 Introduction

Spreading of antibiotic-resistant bacterial pathogens is ever continuing, and the absence of new antibiotics in development pipelines is seriously threatening the future of public health. Vancomycin is a widely used glycopeptide antibiotic for the treatment of infections caused by multi-drug resistant Gram-positive pathogenic bacteria. However, resistance emerged in *Enterococcus* species and now spreads to other bacteria including *Staphylococcus aureus*, causing serious problems in the clinic.<sup>1</sup>

Vancomycin acts by inhibiting peptidoglycan synthesis. The antibiotic interacts with the D-Ala-D-Ala *terminus* of N-acetyl-muramyl-L-Ala-D-Glu-L-Lys-D-Ala-D-Ala late peptidoglycan precursors, hence sequestering the D-Ala-D-Ala dipeptide, and inhibiting the activity of the transpeptidases.<sup>2</sup> Resistance to vancomycin results mainly from the production of modified precursors ending with D-Ala-D-Lac, which exhibits 1000 fold lower binding affinities to vancomycin than D-Ala-D-Ala precursors (Figure 1). Synthesis of D-Ala-D-Lac requires the presence of a ligase with an altered specificity (VanA)<sup>3</sup> that acts at a critical step, thus, reprogramming peptidoglycan synthesis. As a result, it appears as a target of choice to develop new antibiotics. Inhibitors have been discovered<sup>4</sup> on the related enzyme, the D-Ala-D-Ala ligase.

The X-ray crystallographic structure of the D-Ala:D-Lac ligase, VanA from *Enterococcus faecium* (PDB entry: 1E4E) (Figure 2), and that of the D-Ala:D-Ala ligase TtDdl from *Thermus thermophilus* (PDB entry: 2YZG) display similar features. These enzymes are divided in three domains: N-terminal ([A2-G121] and [M1-G104], in blue), central ([C122-S211] and [A105-L192] in red and yellow) and C-terminal ([G212-A342] [S193-T319] in black and green), respectively, in 1E4E and in 2YZG structures. The  $\omega$ -loop (in green in Figure 2) is part of the C-terminal domain. It encompasses the residues

67 [L236-A256] in 1E4E and the residues [Y218-A234] in 2YZG.<sup>6-8</sup> The region opposite to  
68 the  $\omega$ -loop in the structure (yellow in Figure 2) is called “opposite domain” in the present  
69 work. It is composed of residues [A149-Q208] in 1E4E and [V131-K190] in 2YZG and  
70 folds in a two layer  $\beta$  sandwich. The substrates bind to a large pocket located at the  
71 interface between N-terminal, central, and C-terminal domains. In 1E4E, the  $\omega$  loop  
72 closes the pocket and prevents ATP hydrolysis. Conversely in TtDdl, this flexible loop  
73 displayed various extensions in structures obtained with different reaction intermediate  
74 co-crystals (PDB entries: 2YZG, 2YZN, 2ZDG, 2ZDH, 2ZDQ, 2YZM). Cysteines 52 and  
75 64 form a disulfide bridge in crystal structure 1E4E (Figure 2). The bridged form is  
76 called VanA<sub>SS</sub> in the current work.<sup>6,8</sup>

77 The conformational transition of the  $\omega$  loop, inferred from the different ligase struc-  
78 tures, is expected to play a key role in substrates binding, and should thus be studied  
79 in view of a VanA inhibitors development. Conformational transitions of biomolecules  
80 have been extensively studied by molecular modeling,<sup>?</sup> but the impact of these transi-  
81 tions on ligand docking have been investigated less systematically. Here, we propose to  
82 use an Artificial Neural Network, the self-organizing maps (SOMs),<sup>?</sup> to simultaneously  
83 characterize conformational transitions and ligand interactions. SOMs have been used in  
84 the past for the *in silico* screening of chemical compounds for drug discovery,<sup>?,?</sup> for the  
85 prediction of compounds selectivity,<sup>?</sup> for the detection of new bioactive molecules,<sup>?,?</sup> for  
86 the re-scoring of docking poses,<sup>?,?</sup> and for various clustering of conformational ensem-  
87 bles,<sup>?,?,?,37</sup> and protein fragments.<sup>?</sup> A detailed overview of the use of self-organizing maps  
88 in the framework of molecular modeling and structure-based drug design, has recently  
89 been published.<sup>?</sup>

90 The purpose of the present work is: (i) to describe the first steps of the  $\omega$ -loop opening

91 through the analysis of the protein internal dynamics and (ii) to correlate the conforma-  
92 tions sampled along this transition with the binding of ligands displaying various bio-  
93 logical functions. The protein conformational transitions were analyzed with molecular  
94 dynamics simulations, while ligand binding was investigated by molecular docking calcu-  
95 lations. The clustering methods, based on self-organizing maps (SOM) were developed  
96 to cluster protein conformations as well as to classify the ligand poses.

## 97 **3 Materials and Methods**

### 98 **3.1 Preparation of simulated systems**

99 All systems (see Table 1) were setup from the PDB X-ray crystallographic structures  
100 2YZG, 2ZDH and 1E4E corresponding respectively to:

- 101 (i) the D-Ala:D-Ala ligase apo from *Thermus thermophilus* HB8 with open  $\omega$ -loop,
- 102 (ii) the D-Ala:D-Ala ligase from *Thermus thermophilus* HB8 with closed  $\omega$ -loop with  
103 ADP and D-Ala in its binding pocket<sup>6</sup> and,
- 104 (iii) VanA, the D-Ala:D-Lac ligase from *Enterococcus faecium* BM4147, containing ADP  
105 and phosphinate (1(S)-aminoethyl-(2-carboxypropyl)phosphoryl-phosphinic acid).<sup>8</sup>

106 The PDB structure 1E4E was used to produce the systems VanA<sub>SS</sub>.lig bearing an ADP,  
107 a phosphinate inhibitor (PHY) and two Mg<sup>+2</sup> ions in the catalytic site, and a C52-C64  
108 disulfide bridge. The ligands were removed from 1E4E to build the corresponding apo  
109 system, VanA<sub>SS</sub>. Then, the VanA<sub>SS</sub> disulfide bridge was reduced to build the VanA  
110 system. Similarly, the TtDdl<sub>closed</sub>.lig and the TtDdl<sub>closed</sub> systems were built from the  
111 2ZDH structure with or without the ADP, D-Ala and Mg<sup>+2</sup> ions, respectively. Finally,  
112 the TtDdl<sub>open</sub> system was built from the 2YZG structure.

113 Hydrogen atoms were added with the LEaP<sup>9</sup> module of AMBER 10.<sup>10</sup> The FF99SB  
114 force field<sup>11</sup> was used. The systems were neutralized with Na<sup>+</sup> counter-ions. The organic  
115 molecules were parametrized with Antechamber<sup>12</sup> and the **General AMBER Force Field**  
116 (**GAFF**).<sup>13</sup> Explicit TIP3P<sup>14</sup> solvent water molecules were added to the systems in a  
117 cubic box under periodic boundary conditions with a buffer zone of 10Å. The system  
118 components are given in Table 1.

## 119 **3.2 Molecular Dynamics Simulations**

120 The **Simulated Annealing with NMR-Derived Energy Restraints** (Sander) module from  
121 AMBER 10<sup>15</sup> was used to perform five rounds of minimizations composed of steepest  
122 descent followed by conjugate gradient algorithms. Harmonic restraints were applied on  
123 the protein atom position with the reference set to the final position of the previous  
124 round and a force constant of 100, 50, 25, 10 and 5 kcal · mol<sup>-1</sup> · Å<sup>-2</sup> in each round,  
125 respectively. Then, the systems were thermalized to 298 K for 20 ps with **Molecular**  
126 **Dynamics** (MD) at constant volume, by making use of the weak-coupling algorithm<sup>16</sup>  
127 and harmonic restraints of 25 kcal · mol<sup>-1</sup> · Å<sup>-2</sup> on the solute atom positions. Thus, six to  
128 seven equilibration rounds were performed with a Langevin thermostat with a collision  
129 frequency  $\gamma = 2 \text{ ps}^{-1}$ . One 5 ps MD round at constant volume was followed by four  
130 2.5 ps and one 10 ps constant pressure MD rounds. Harmonic restraint force constants  
131 were 25, 25, 20, 15, 5 and 2.5 kcal · mol<sup>-1</sup> · Å<sup>-2</sup>, respectively. Finally a last MD round of  
132 60 ps was performed without any restraints.

133 **Molecular Dynamics** (MD) trajectories were recorded over 20 to 30 ns with the  
134 **Particle Mesh Ewald Molecular Dynamics** (PMEMD)<sup>?,17</sup> module from AMBER 10.  
135 A cutoff of 10 Å was used for Lennard-Jones interaction calculations. Long-range elec-

136 trostatic interactions were calculated with the **P**article **M**esh **E**wald (PME) protocol.<sup>17</sup>  
137 The simulations were performed at a pressure of 1 atm and a temperature of 298 K under  
138 the control of a Berendsen thermostat with a coupling time of 2 ps.<sup>16</sup> The SHAKE algo-  
139 rithm<sup>18</sup> kept all covalent bonds involving hydrogens rigid so integration time step of 2 fs  
140 was used for all MD simulations. Atomic coordinates were saved every picosecond. D-  
141 Ala:D-Lac ligase MD trajectories were recorded seven to nine times with different initial  
142 random seeds. The D-Ala:D-Ala ligase trajectories were recorded only once.

### 143 **3.3 Conformational analysis of the molecular dynamic simula-** 144 **tions using self-organizing maps**

145 **S**elf-**O**rganizing **M**aps (SOM),<sup>19,20</sup> which are unsupervised neural networks, were used  
146 to cluster the 50 000 conformations sampled during the “VanA” and “VanA<sub>SS</sub>” MD sim-  
147 ulations. Conformations were encoded as follow: the  $n \times n$  pairwise square Euclidean  
148 distance matrix  $D$  was calculated for  $n$   $C_\alpha$  atoms of the protein. Then, to compress the  
149 information, the covariance matrix,  $C$  of the lines versus columns of  $D$  was calculated:<sup>21</sup>

$$C_{i,j} = \frac{1}{n} \sum_{k=1}^n \sum_{l=1}^n (d_{i,k} - \bar{d}_i)(d_{l,j} - \bar{d}_j) \quad (1)$$

150 where  $\bar{d}_i = \frac{1}{n} \sum_{j=1}^n d_{i,j}$ . As  $C$  describes a 3D object, its eigenvalues beyond the first  
151 four are null. Hence, the eigenvectors of  $C$ ,  $N_{i=1,\dots,4}$ , corresponding to the four first  
152 eigenvalues, were kept applied to  $D$ ;  $D \cdot N_{i=1,\dots,4}$ . This compression in  $n \times 4$  matrices gives  
153 a conformational descriptor, which conserves information.

154 These descriptors were used to train a periodic Euclidean self-organizing map (SOM).  
155 Most commonly used SOMs are 2D SOMs, which are defined by three-dimensional ma-

156 trices. The first two dimensions, 2D, lengths are chosen by the user, here  $50 \times 50$ , and  
157 define the map size. As these dimensions are chosen to be periodic, the map is a toroid.  
158 The third dimension has the length of the input vectors, or descriptor, here:  $4n$ , and each  
159 vector along the third dimension is called a neuron.

160 The self-organizing maps were initialized with a random uniform distribution covering  
161 the range of values of the input vectors. At each step, an input vector is presented to the  
162 map, and the neuron closest to this input, the Best Matching Unit (BMU) is updated.  
163 The maps were trained in two phases. *Guillaume: peux tu revoir cela: During the*  
164 *first phase, the 50 000 input vectors are presented to the SOM in random order to avoid*  
165 *mapping bias with a learning parameter of 0.5, and a radius parameter of 36, as explained*  
166 *in Reference 40.*

167 *During the second phase, the learning and radius constants were decreased exponen-*  
168 *tially from starting values 0.5 and 36, respectively, during 10 cycles of presentation of all*  
169 *the data in random order.*

170 Hence, to delineate clusters on the SOMs, the conventional **Unified distance matrix**  
171 (U-matrix) is a useful tool. For each neuron  $\nu$  on the map, a corresponding U-matrix  
172 element is calculated as the mean Euclidean distance between the neuron  $\nu$  and its eight  
173 immediate neighbors:

$$\text{U-height}(\nu) = \frac{1}{8} \sum_{\mu \in N(\nu)} d(\nu, \mu) \quad (2)$$

174 where  $N(\nu)$  is the set of neighbors, and  $d(\nu, \mu)$  is the Euclidean distance between the  
175 vectors  $\mu$  and  $\nu$ . The resulting  $50 \times 50$  U-matrix reveals the topological organization  
176 of the map, and can be used to draw the contours of clusters by applying a threshold  
177 distance value.

178 SOMs distribute data on the map so that points which are close or far in the descriptor

179 space are also close or far, respectively, on the map. However, they also distribute the data  
180 as evenly as possible on the map. This action enforces the similarity between neighboring  
181 neurons in the final map. If the system topology is poorly compatible with a projection  
182 on a torus, some distant conformational basins will be projected on close regions of the  
183 SOM, resulting in large conformational variations between close neurons, which, coupled  
184 with the enforcement of similarity just described, induce the formation of empty nodes.

185 The maps convergence was assessed quantitatively, by running 80 independent SOM  
186 calculations. Each calculation started from a different random map, and the comparison  
187 of the resulting maps was performed using the following flooding algorithm, inspired by  
188 the watershed algorithm<sup>7</sup> used in image processing. This algorithm works on the topology  
189 of the U-matrix. It starts from the global minimum and flood the map according to the  
190 landscape of the U-matrix. The maps are then reordered according to the order of the  
191 flooding process. The maps were compared by calculation of the average correlation  
192 between the reordered neurons. Since these averages were within the interval 0.98-1.0  
193 (data not shown OR FIGURE corrdist.pdf), the map convergence was considered as  
194 effective, and the maps valid.

195 Representative conformations extracted from the SOMs clusters are available from  
196 the authors upon request.

### 197 **3.4 Flow analysis of the self-organizing maps**

198 The molecular dynamic trajectory evolution can be followed for each time step  $t$  by its  
199 position  $(i,j) = \Phi(t)$  on the SOM. The ensemble time steps that project on neuron  $(i,j)$   
200 is called  $\{\tau_{i,j}\}$ , and the total number of these steps is noted  $f_{i,j}$ . The local mean transfer  
201 vector field is then defined by the average SOM index difference for times  $t$  in  $\tau_{i,j}$  to the

202 next steps  $t + 1$ :

$$\mathbf{v}_{i,j} = \frac{1}{f_{i,j}} \sum_{t \in \tau_{i,j}} \frac{\Phi_{t+1} - \Phi_t}{\|\Phi_{t+1} - \Phi_t\|} \quad (3)$$

203 only defined for non-empty neurons where  $f_{i,j}$  is non zero.

### 204 3.5 3D Self Organizing Maps

205 Similarly to the 2D SOM described in the previous sections, 3D self-organizing maps were  
206 built to describe the docking position of the ligands atoms.

207 The input of that SOM procedure was the set of 3D coordinates for individual atoms  
208 of the ligand along the molecular dynamics trajectory. The 1170000, 121199 and 275000  
209 input vectors for ADP, D-Ala and D-Lac, respectively, were used to train three inde-  
210 pendent 3D SOMs. The self-organizing maps were initialized with a random uniform  
211 distribution of ligand coordinates and trained in two phases. During the first phase,  
212 input vectors are presented to the SOM in random order. **Guillaume, peux tu donner**  
213 **le taux de décroissance... Initial radius and learning parameters were set to 7.5 and 1.0,**  
214 **respectively, and decreased exponentially to 0 during the training process.** As described  
215 before, the SOM convergence was checked by multiple independent training runs.

216 Different sizes (ADP: 25x17x18; D-Ala: 18x17x11; D-Lac: 17x15x13) were chosen  
217 for the 3D SOMs, corresponding respectively to resolutions of 10.28, 60.98 and 38.51  
218 neurons/ $\text{\AA}^3$ , and thus to: 2.2, 3.9 and 3.4 neurons/ $\text{\AA}$ . Such resolutions in neurons cor-  
219 respond to a precision of 0.25-0.5  $\text{\AA}$  in atomic coordinates, similar to the estimated  
220 positional error in X-ray crystallographic structures at about 2.5  $\text{\AA}$  resolution.

## 221 **3.6 Docking procedure**

222 The ATP, D-Ala, D-Lac, D-alanylphosphate (D-Ala(P)), the phosphinate (PHY, tran-  
223 sition state inhibitor) and D-Ala-D-Lac, which are involved in, or interfere with VanA  
224 enzymatic activity, were formatted in mol2 with Chimera 1.4<sup>22</sup> and MarvinSketch 5.1<sup>23</sup>  
225 for docking.

226 UCSF DOCK 6.5<sup>24-26</sup> was used to perform ligand docking on representative VanA<sub>SS</sub>  
227 MD conformations selected by 2D SOM analysis. These structures were those having  
228 their structure descriptor closest in Euclidean distance from that of a populated neuron.  
229 Chimera<sup>22</sup> was used to add hydrogens, check atom assignment, and assign partial charges  
230 in line with the AMBER-ff99SB force field. It was also used to produce mol2 format  
231 files for the ligands and the selected conformations of the receptor. The DMS software  
232 program<sup>27,28</sup> generated the molecular surface of the receptor using a radius probe of 1.4 Å.  
233 Then, spheres were calculated around the receptor with the DOCK 6.5 command 'sphgen'  
234 with radius probe values varying between 1.4 Å and 4 Å. Spheres within a radius of 10 Å  
235 around the geometric center of the crystallographic ligands (ADP, PHY) found in 1E4E  
236 were selected. The grid encoding van der Waals and electrostatic interactions was pre-  
237 calculated with the "grid" tool<sup>29</sup> in a box containing the selected spheres. The DOCK  
238 program builds up to 500 flexible ligand docking poses, on the pre-calculated "grid"  
239 interaction map. The ligand poses were then re-scored with the implementation of the  
240 Hawkins **M**olecular **M**echanics **G**eneralized **B**orn **S**urface **A**rea (MM-GBSA) score,<sup>30-34</sup>  
241 implemented in UCSF DOCK 6.5.

242 The best scoring solution was kept for each protein - ligand pair. The binding pocket  
243 was defined by residues: E14, E15, V18, H98, G99, E103, S126, C129, M130, K132, T135,  
244 Y136, K170, P171, S174, G175, S176, S177, F178, V180, E213, I239, F240, R241, I242,

245 H243, Q244, R289, D291, L302, N303, E304, V305, N306, T307, P309, G310, S315, R316  
246 and Y317.

## 247 4 Results

### 248 4.1 Concerted $\omega$ -loop / opposite domain motions correlate with 249 the presence of the disulfide bridge

250 The global **R**oot **M**ean **S**quare **D**eviation (RMSD) from the initial structure for the  $C\alpha$   
251 atoms stabilized at about 2.2 Å for the eight independent VanA trajectories (Figure 3a)  
252 and the ten VanA<sub>SS</sub>.lig MDs (Figure 3c). By contrast, the seven VanA<sub>SS</sub> MDs (Figure  
253 3b) displayed heterogeneous behavior. Curve with the smallest drift for this system, in  
254 black, is similar to that observed for VanA, whereas that with the largest drift, in red,  
255 increased up to 3.5Å after 17 ns (Figure 3b). Hence, the presence of the C52-C64 disulfide  
256 bridge correlated with a destabilization of VanA conformations.

257 The contributions of the different regions (C-terminal, central, N-terminal, opposite  
258 domains and  $\omega$ -loop) to the RMSD were analyzed on the trajectories with the largest  
259 global RMSD drifts recorded for VanA, VanA<sub>SS</sub> and VanA<sub>SS</sub>.lig (Figure 3d-f). A similar  
260 analysis was performed for the D-Ala:D-Ala ligase systems TtDdl<sub>open</sub>, TtDdl<sub>closed</sub> and  
261 TtDdl<sub>closed</sub>.lig (Figure 3g-i, Table 1). The  $\omega$ -loop always displayed the largest drift,  
262 except for TtDdl<sub>closed</sub>.lig (Figure 3i). The systems displaying the smallest  $\omega$ -loop drifts  
263 were VanA and TtDdl<sub>closed</sub> (Figure 3d & h).

264 The large drifts of the  $\omega$ -loop in MD simulations are in good agreement with the large  
265 conformation differences observed in X-ray structures.<sup>6</sup> Indeed, the  $\omega$ -loop covers the  
266 binding site entrance in 2ZDH and 1E4E, whereas it extends away from the core of the D-

267 Ala:D-Ala ligase in 2YZG. Correspondingly, the largest observed drift was for TtDdl<sub>open</sub>,  
268 (2YZG), which also has an empty catalytic site, and is probably in an inactive functional  
269 state. Interestingly, among the three systems built from 1E4E, VanA<sub>SS</sub> (Figure 3e) and  
270 VanA<sub>SS</sub>.lig (Figure 3f) presented large  $\omega$ -loop drifts despite an initial closed conformation.  
271 Noticeably, the large global protein RMSD drift observed in the presence of the C52-C64  
272 disulfide bridge (Figure 3b,c), is mostly due to the  $\omega$ -loop motions (Figure 3e,f). In the  
273 presence of the substrates ADP and PHY, the disulfide bridge still destabilized the  $\omega$ -  
274 loop, but to a lesser extent (Figure 3c,f). However, the presence of the ligands strongly  
275 reduced all protein region drifts when the  $\omega$ -loop is wrapped (Figure 3i).

276 The other protein regions rarely drifted beyond 3Å. Nonetheless, the opposite domain  
277 (yellow curves), the central domain (red curves) and the C-terminal domain (black curves)  
278 drifted more when the  $\omega$ -loop made large motions (Figures 3e-g).

279 To describe the relative displacement of the protein regions with respect to each other,  
280 a **Principal Component Analysis** (PCA) was performed on the C $\alpha$  atoms trajectories  
281 (Figure 4). More eigenvectors were necessary to account for 90 % of the motions of the  
282 “opened” systems, than for the “closed” systems, with 16 to 25 and 46 to 61 eigenvectors  
283 required, respectively. A large and strongly dominant eigenvalue was observed (Figure 4f)  
284 for the simulations VanA<sub>SS</sub> and VanA<sub>SS</sub>.lig, which displayed strongly correlated motions,  
285 and large  $\omega$ -loop drifts. The relative importance of the first eigenvalue was lower in  
286 the presence of ligands. For example, the first eigenvector of VanA<sub>SS</sub> and VanA<sub>SS</sub>.lig,  
287 contributed respectively 18.4 % and 6.0 % to the global motion variance. For VanA  
288 and TtDdl<sub>closed</sub> simulations, no dominant motion was observed as the first eigenvalue  
289 accounted for 1.5 % to 3.0% of the global motion.

290 The projection of the first PCA mode on the protein structures (Figures 4a-e) showed

291 homogeneously distributed motions with relatively small amplitude in TtDdl<sub>closed</sub> and  
292 VanA closed state MD simulations (Figures 4c,d). By contrast, motions were mainly  
293 located in the  $\omega$ -loop and the opposite domain for VanA<sub>SS</sub>, VanA<sub>SS</sub>.lig and TtDdl<sub>open</sub>,  
294 (Figures 4a,b,e).

295 Hence, PCA analysis revealed the specific internal fluctuations of the  $\omega$  loop and the  
296 opposite domain. As expected, these fluctuations are larger for structures bearing an  
297 open  $\omega$  loop and no ligand in the catalytic site. It was more surprising to find that the  
298 C52-C64 disulfide bridge would also increase so significantly the  $\omega$  motions.

## 299 4.2 Self-organizing maps suggest contours of free-energy basins

300 Self-Organizing Maps (SOM)<sup>19,20</sup> were used to project the conformational space explored  
301 by the ligase during MD simulations onto a smaller, bi-dimensional and topologically  
302 organized space. A  $50 \times 50$  SOM was trained to cluster the protein conformations along  
303 one trajectory of VanA and VanA<sub>SS</sub>, respectively. The VanA and VanA<sub>SS</sub> MD trajectories  
304 analyzed here occupied distinct zones of the SOM.

305 As described in Materials and Methods, the U-matrix is a convenient visualization  
306 tool to reveal SOM topological features.<sup>35,36</sup> Closely related structures are grouped in  
307 the same valleys or basins with small inter-neuron distances colored in blue separated by  
308 ridges of large inter-neuron distances, defining their boundaries in red (see Figure 5a).  
309 The U-matrix, which gives an evaluation of the state density, can thus be interpreted  
310 as a qualitative marker of the free energy landscape of the protein conformational space  
311 within the sampled area. The landscapes of VanA and VanA<sub>SS</sub> showed large blue patches  
312 of homogeneous structures separated by thin red barriers that would be expensive to  
313 cross and lower green walls that can be crossed occasionally. For VanA, which performed

314 limited exploration, there were few big clusters separated by low barriers (called 4 in  
315 Figure 5a). For VanA<sub>SS</sub>, there were larger barriers roughly dividing the U-matrix into  
316 two main regions, the first bearing two sub-regions (depicted by 1 and 2 in Figure 5a)  
317 and the second displaying a higher degree of diversity (noted 3 in Figure 5a). The map  
318 showed that the VanA MD spanned a smaller physical space (mostly blue neurons) than  
319 VanA<sub>SS</sub> which formed at least two independent coherent tracts (basins 1,2 *versus* 3) with  
320 higher diversity in the second one according to intrinsic distance (cyan to green neurons).

321 To give a quantitative support to the interpretation of the SOM clusters as free-  
322 energy basins, MM-PBSA and MM-GBSA energies were calculated along the VanA and  
323 VanA<sub>SS</sub> trajectories, with the AMBER 12 package tools,<sup>?</sup> and projected on the SOM  
324 (Figure 7a,b). Although MM-PBSA and MM-GBSA energies displayed significant fluc-  
325 tuations, they agreed reasonably well with the SOM clustering, since energies were more  
326 uniform within basins than between them and the contiguous basins borders displayed  
327 higher energies. This relative correspondance between SOM clustering and the energies  
328 supports the U-matrix as qualitative marker of the free energy landscape.

329 Structural properties were then projected and visualized onto the 2D trained map.  
330 The projection of the RMSD from the first frame of each trajectory further corroborates  
331 the quality and the convergence of the clustering process (Figure 5b) and the relation  
332 with the conformational landscape. For VanA<sub>SS</sub> the two regions delineated by the U-  
333 Matrix displayed distinct RMSD values. The lower U-matrix zone (basins 1 and 2)  
334 corresponded to comparable drifts to that of VanA, while the higher U-matrix zone (basin  
335 3) revealed conformations that had largely evolved from the initial structure (Figure 5b).  
336 The homogeneous RMSD pattern of VanA and the bipolar one of VanA<sub>SS</sub> are directly  
337 related to the U-matrix patterns.

338 In order to evaluate structural changes, the evolution of the  $\beta$ -strands secondary struc-  
339 ture content was projected on the SOM (Figure 5c). The  $\alpha$  helices were only marginally  
340 affected on both systems. Only a slight uncoiling was observed for VanA between 10 and  
341 15 ns of simulation (data not shown). Up to 12 amino acids loose their  $\beta$  structure in  
342 VanA<sub>SS</sub> and VanA trajectories as can be seen on the SOM projection (Figure 5c). The  
343  $\beta$ -6 strand located in the N-terminal domain of the protein (Figure 2) lost three to four  
344 residues in VanA<sub>SS</sub> and no more than three in the last part of VanA trajectory. The  
345 most affected  $\beta$  structures apart from  $\beta$ -6 were  $\beta$ -14 and  $\beta$ -15, close to the  $\omega$ -loop  
346 (Figure 2). Indeed, while  $\beta$ -15 gained 3 residues in VanA, both  $\beta$ -14 and  $\beta$ -15 lost two  
347  $\beta$  residues in the most-drifting part, the last 10 ns, of the MD trajectory of VanA<sub>SS</sub>. This  
348 secondary structure variability agrees with the role of hinges played by these  $\beta$  strands  
349 during the opening motion of the  $\omega$ -loop in VanA<sub>SS</sub>, as can be seen in (Figures 4a,c).

350 The SOM appeared to produce a meaningful clustering of conformations. Since the  
351 original trajectory can be followed on the map, SOM could also be used to investigate  
352 how the protein evolves in the different parts of the map with the vectors field  $\mathbf{v}_{i,j}$  (see  
353 Materials & Methods Eq. 3). The vector field gives the propensity of the mapped  
354 conformations to evolve in the given direction. The vectors field appeared to follow  
355 the gradient of the U-matrix. The vectors with low or null norms are mostly present  
356 in the bottom of the basins. These results substantiate the interpretation of the U-  
357 matrix as a marker of the free energy landscape. Large arrows indicate high net flow for  
358 densely populated regions, but could also be due to poor statistics in low density regions.  
359 Interestingly, some small vectors are also present on the U-matrix barriers between the  
360 closed and opened conformations of VanA<sub>SS</sub>. These structures have the same probability  
361 to go to either basins, which, in practice, would correspond to the definition of a transition

362 states ensembles.<sup>38</sup> The trajectory of VanA<sub>SS</sub> (in pink in Figure 6) is characterized by  
363 two major basins. The first basin, subdivided in two sub-basins, 1 and 2, groups initial  
364 and then more equilibrated conformations of the closed state of VanA<sub>SS</sub> respectively. The  
365 second major basin, labeled “3”, contains open states. The barrier between basins 2 and 3  
366 is composed of low density neurons, with high convergent flows pointing to the transition  
367 states ensemble surrounded by stationary points. As already seen, VanA and VanA<sub>SS</sub>  
368 covered distinct conformational spaces except for a limited border area highlighted by  
369 brown diamonds in Figure 6.

370 The transition states between basins are defined as points of zero flow. Interestingly  
371 such points correspond to saddle point in the surface of the U-matrix. As explained in  
372 the Materials and Methods, flow is not defined at empty neurons, and thus points of zero  
373 flow close to empty neurons should not be picked up as saddle points. Saddle points  
374 detected in the present work (Figure 6) are located far from empty neurons.

375 The conformational clustering of the molecular dynamics simulations VanA and VanA<sub>SS</sub>  
376 show several basins corresponding to closed and open conformations of VanA. The de-  
377 tection of such conformations gives a more precise picture of the different steps of the  
378 interaction between VanA and the reaction substrates and should help to search for VanA  
379 inhibitors.

380 Hence, density metrics given by the U-matrix suggests that the basins could be inter-  
381 preted or defined as free-energy basins, within the limits of the conformational sampling.  
382 Projection of  $\beta$  secondary structure evolution indicated that the  $\beta$  strands located close  
383 to the  $\omega$ -loop hinges were the most variable ones. The analysis of the conformational flow  
384 defined populated regions during the  $\omega$  loop opening that can be considered as transition  
385 states.

386 As described by 6 for TtDdl the opposite domain of VanA<sub>SS</sub> moves away from the  
387 binding cavity (Figure 5d, 4a). In contrast to the observed motions in TtDdl, only a sub-  
388 part of VanA<sub>SS</sub> central domain, the opposite domain, is involved in the opening motion  
389 during the course of the dynamics.

### 390 **4.3 SOM classification of ligand poses related to their function**

391 In the Ter-Ter mechanism of the ligases (Figure 1),<sup>?</sup> the ATP binds first. It is followed by  
392 a first D-Ala and then either D-Lac, or a second D-Ala for VanA or Ddl, respectively. The  
393 ligands (ATP, D-Ala, D-Lac, PHY, D-Ala(P), D-Ala-D-Lac) were docked individually on  
394 conformations representing each neuron of the 2D-SOM to relate conformations sampled  
395 along the  $\omega$  loop opening and ligand binding propensity. A neuron was represented by  
396 the structure, either from VanA or VanA<sub>SS</sub>, which had closest descriptor to that of the  
397 neuron after training.

398 One 3D self-organizing map, 3D-SOM, was built from the docking results for each  
399 ligand. The descriptors were the coordinates of all atoms of each ligand. Mapping the  
400 identity of the ligand, (ADP, D-Ala, *etc.*) on the resulting map indicated their respective  
401 consensus binding sites. The 3D-SOM was projected onto the 3D Cartesian coordinates  
402 simply using the neuron descriptor field (Figure 8). The respective ligand binding sites  
403 agreed with those observed in crystal structures of TtDdl in complex with ADP and D-Ala  
404 (2ZDH).<sup>6</sup> Interestingly, the binding sites identified by docking here for D-Lac overlapped  
405 with those of phosphinate, a transition state analog co-crystallized with VanA (1E4E) or  
406 that of D-Ala-D-Ala in TtDdl (2ZDQ).<sup>6,8</sup>

407 To further analyze the ligand docking specificity, the GB/SA docking scores (see Ma-  
408 terials and Methods) were then projected onto the 2D SOM used to cluster the MD

409 conformations (Figure 9). Noticeably, VanA displayed binding trends that agreed with  
410 the enzymatic role of the ligand. In addition to the ligand binding site specificity ob-  
411 served with the 3D-SOM, the conformational ligand binding specificity could hence be  
412 established. For instance, ATP binds exclusively in the opened  $\omega$ -loop conformation basin  
413 defined by the U-matrix (Figures 5 & 9a) and scored better than the reaction products  
414 ADP (data not shown). The second partner of the reaction, D-Ala, binds non-selectively  
415 to almost all the VanA<sub>SS</sub> structures (Figure 9b), and less than half of the VanA struc-  
416 tures. Not surprisingly, the product of the enzymatic reaction, D-Ala-D-Lac (Figure 9f),  
417 does not display binding selectivity. The acylphosphate, D-Ala(P) (Figure 1) correspond-  
418 ing to the phosphorylated form of the former D-Ala, and D-Lac only binds with a good  
419 score, to the same restricted region of the SOM map, corresponding to the third basin  
420 where the reaction takes place (see Figure 9c,d). Phosphinate mimicking the tetrahedral  
421 intermediate binds also with a good score to the third basin (see Figure 9e). Strikingly,  
422 the best phosphinate docking scores were observed on the conformations, that were delin-  
423 eated as the transition state ensemble between closed and open  $\omega$ -loop states in section  
424 “Self-organizing maps suggest contours of free-energy basins”.

## 425 5 Discussion

426 In the present work, we used MD simulations to investigate VanA conformational sam-  
427 pling, in particular the first opening steps of the  $\omega$ -loop. Two main conformational basins  
428 were visited in the presence of a disulfide bridge between C52 and C64. Known ligands  
429 (substrates, products and intermediate alike) were docked on representative conforma-  
430 tions issued from the clustering. This analysis showed a correlation between docked  
431 ligand binding energies and the protein conformation, which is in good agreement with

432 the Ter-Ter ordered mechanism of the ligase.

433 The MD simulations performed here indicated that the  $\omega$ -loop opening mechanism  
434 of VanA is similar to that of the endogenous enzyme, TtDdl.<sup>6,45</sup> Indeed, the semi-open  
435  $\omega$ -loop conformation of VanA is similar to that observed in the 2ZDG TtDdl structure.<sup>6</sup>  
436 Furthermore, the correlation of the  $\omega$ -loop and opposite-domain motions (Figure 4) in the  
437 MD simulations agreed with available structures data on ligases.<sup>6</sup> The similarity between  
438 consensus binding sites of ADP, D-Ala and D-Lac in representative conformations of  
439 VanA (Figure 8) and those observed in crystal structures of TtDdl (2ZDG, 2ZDH, 2ZDQ)<sup>6</sup>  
440 strongly supports that those two proteins make similar interactions with their substrates.  
441 This ligand binding similarity and the mechanistic similarities implied by MD simulations  
442 interestingly supports the idea that new inhibitors against both D-Ala:D-Ala and D-  
443 Ala:D-Lac ligase could be found and developed.

444 A recent study of the D-Ala:D-Ala ligase described a possible  $\omega$ -loop opening mech-  
445 anism in Ddl by **Steered Molecular Dynamics** (SMD).<sup>45</sup> Conformations extracted from  
446 this opening path were used in an initial screening, which allowed to identify experimen-  
447 tally validated inhibitors. This study highlighted the importance of the  $\omega$ -loop opening  
448 conformational analysis in the quest for new ligase inhibitors. In addition, the impor-  
449 tance of the  $\omega$ -loop dynamics for the D-Ala:D-Lac ligase, was shown. Furthermore, the  
450 opposite domain motion is also crucial for the activity of the VanA ligase.

451 The clustering of molecular dynamics simulations, performed here using SOMs, was  
452 used to extract representative conformations. The representative conformations have  
453 different propensities to bind ligands at different stages of the enzymatic reaction (sub-  
454 strates, intermediate-like, products), as it was shown by the clustering of the ligand  
455 docking poses. These conformations are thus good candidates to perform virtual screen-

456 ing runs in the context of the development of new antibiotics able to overcome pathogenic  
457 resistance.

458 The new insights into the relationship between VanA conformational transition and  
459 predicted ligand interactions were made possible by the use of the self-organizing maps  
460 (SOMs).

461 The main advantage of the distance matrix based SOM, compared the usage of Carte-  
462 sian coordinates,<sup>37</sup> is that the clustering is independent of any structural alignment. This  
463 is of major importance to cluster structures involving large conformational changes as in  
464 protein folding studies.<sup>40</sup>

465 However, distance matrices are highly redundant, and PCA compression<sup>21</sup> was used  
466 to reduce data size. Finally, the SOM algorithm, applied to PCA compressed distance  
467 matrices,<sup>19,20</sup> gives rise to a conformational clustering method that is independent of any  
468 choice of reference conformation, or any coordinate RMSD calculation.

469 Another advantage of self-organizing maps is that they provide a simplified description  
470 of the conformational space of a protein, without having to choose specific variables  
471 describing the principal motions.

472 However, a limitation in the interpretation of the U-matrix in terms of free energy  
473 landscape and transition state ensembles arises in the present study from the length of the  
474 molecular dynamics trajectories. 25 ns is a short time interval compared to the timescales  
475 usually simulated when one performs a full analysis of the free energy landscape for the  
476 system.<sup>?</sup> A quantitative analysis of the convergence of the trajectories in each basin  
477 determined from the SOM clustering was attempted by using the cosine content.<sup>?</sup> Values  
478 of 0.105, 0.813 and 0.929 are respectively obtained on the trajectory VanA and on the two  
479 time intervals of 0-16.2 and 16.2-25 ns of the trajectory VanA<sub>SS</sub> (Figure 3b), before and

480 after opening of the loop  $\omega$ . The small value obtained on the trajectory VanA agrees with  
481 the short timescale of the oscillatory motion observed for the  $\omega$  loop in this trajectory.  
482 In contrast, along VanA<sub>SS</sub>, more complex dynamical behavior is observed, which is not  
483 dominated by one single motion. Because of this complexity, the motion timescales cannot  
484 be efficiently sampled during the short 25 ns trajectories recorded in the present work.  
485 As most of the trajectories are far from being converged, the prediction of free energy  
486 profiles from the conformational clustering by SOMs should thus be considered as being  
487 only qualitatively.

488 Nevertheless, in the particular case studied here, due to the existence of very relevant  
489 and different X-ray crystallographic structures from the TtDdl ligase, it was possible to  
490 obtain interesting insights into the free energy landscape of VanA.

491 The projection of the RMSD onto the SOM (Figure 5b) revealed a description of the  
492 conformational space dividing the set of conformations into distinct basins, in agreement  
493 with the global RMSD observation along MD trajectories (Figure 3). Furthermore, the  
494 transition structures between the basins can be detected by searching saddle points in  
495 the U-matrix.<sup>41</sup> Interestingly, these transition structures are favorable for the docking of  
496 the phosphinate tetrahedral-intermediate analog (Figure 9e).

497 The conformational clustering by SOMs gives a statistical picture of the MD simula-  
498 tion evolution. Through the preservation of the Boltzmann distribution by the molecular  
499 dynamics, the map resulting from the SOM algorithm contains information on the free-  
500 energy surface of the conformational space. One important feature of the SOM is to  
501 preserve the topological organization of the input space: closely related structures of the  
502 input space are grouped together in the SOM output space. Another trend is to distribute  
503 evenly data on the map so that apart from highly favorable or unfavorable zones the neu-

504 ron occupancy is homogeneously distributed. Hence, within the limits of the sampling  
505 completeness, SOMs seemed to provide a relevant delineation of free energy areas. The  
506 length of the simulations (25 ns) proved sufficient to offer significantly different docking  
507 specificities that could reflect the function of the ligands.

508 The relation between SOM clustering and conformation propensities suggests that  
509 SOMs could give a general framework for the definition of relevant reaction coordinates  
510 or collective variables allowing readily to project the evolution of MDs on the free energy  
511 topological map. Hence, SOM clustering appears attractive to analyze the conformational  
512 sampling in the framework of enhanced sampling methods.<sup>42,43</sup> The limits between free-  
513 energy basins are characterized by low populated areas, reflecting a low probability to  
514 access this conformation during the MD simulation. Since SOMs can be used to define  
515 free energy basins they readily allow the identification of the transition state ensembles  
516 by analysis of the flow given by the transfer vectors field, looking for null flow neurons  
517 implying an equiprobability to reach either close-by basin as described by Bolhuis and  
518 Ding<sup>38,41</sup> and Vanden-Eijnden.<sup>44</sup>

519 SOM analysis can also simply relate protein conformation to the ligand binding  
520 propensity by projecting of ligand docking scores on the conformational 2D-SOM. Pose  
521 classifications agreed with the ligand function, which supports the coherence of docking  
522 and scoring. These results validate the docking protocol as a specific tool to identify po-  
523 tential inhibitors of the D-Ala:D-Lac ligase. Furthermore, this analysis allows to choose  
524 the most relevant conformations to search for specific inhibitors by virtual screening. In  
525 the frame of the docking study on D-Ala:D-Lac ligase, taking into account the  $\omega$  loop  
526 flexibility was essential to cluster ligands according to their functions, in agreement with  
527 results recently obtained by molecular docking on D-Ala:D-Lac ligase from *Leuconostoc*

528 *mesenteroides*.?

## 529 **6 Conclusion**

530 Molecular dynamics simulations of the D-Ala:D-Lac ligase was used to investigate the  
531 substrates binding mechanism. First, it appeared that the presence of a disulfide bridge  
532 between cysteines C64 and C52 induced the opening of the  $\omega$ -loop and of the opposite  
533 domain, which is essential for unhindered entrance in the ligase catalytic site. Second, the  
534 development of an original clustering approach delineated the early steps of the open-  
535 ing mechanism and helped to identify representative conformations of this transition.  
536 The docking of known ligands on these representative conformations unraveled the rela-  
537 tion between conformation and docking propensity in agreement with the ligand function.

538

539 We propose self-organizing maps as a general method for relating conformational tran-  
540 sition of biomolecules and ligand docking poses.

541

542 This paves the way for the selection of appropriate binding site and pocket conforma-  
543 tions for the search of D-Ala-D-Lac inhibitors. Furthermore, the conformation clustering  
544 can be related to the definition of system free-energy basins along MD simulations, and  
545 could thus be of interest in the frame of enhanced sampling and conformational free  
546 energy landscape simulations.

## References

- [1] Courvalin, P. Vancomycin resistance in gram-positive cocci. *Clinical Infectious Diseases* **2006**. 42, S25–S34.
- [2] Reynolds, P. E. Structure, biochemistry and mechanism of action of glycopeptide antibiotics. *European Journal of Clinical Microbiology and Infectious Diseases* **1989**. 8, 943–950.
- [3] Arthur, M.; Molinas, C.; Bugg, T.; Wright, G.; Walsh, C.; Courvalin, P. Evidence for in vivo incorporation of d-lactate into peptidoglycan precursors of vancomycin-resistant enterococci. *Antimicrobial agents and chemotherapy* **1992**. 36, 867–869.
- [4] Reynolds, P. E.; Snaith, H. A.; Maguire, A. J.; Dutka-Malen, S.; Courvalin, P. Analysis of peptidoglycan precursors in vancomycin-resistant enterococcus gallinarum bm4174. *Biochemical Journal* **1994**. 301, 5.
- [5] Arthur, M.; Reynolds, P.; Courvalin, P. Glycopeptide resistance in enterococci. *Trends in microbiology* **1996**. 4, 401–407.
- [6] Kitamura, Y.; Ebihara, A.; Agari, Y.; Shinkai, A.; Hirotsu, K.; Kuramitsu, S. Structure of D-alanine-D-alanine ligase from *Thermus thermophilus* HB8: cumulative conformational change and enzyme-ligand interactions. *Acta Cryst D* **2009**. 65, 1098–1106.
- [7] Meziane-Cherif, D.; Saul, F.; Haouz, A.; Courvalin, P. Structural and Functional Characterization of VanG D-Ala: D-Ser Ligase Associated with Vancomycin Resistance in *Enterococcus faecalis*. *J Biol Chem* **2012**. 287, 37583–37592.

- 568 [8] Roper, D.; Huyton, T.; Vagin, A.; Dodson, G. The molecular basis of vancomycin  
569 resistance in clinically relevant enterococci: crystal structure of D-alanyl-D-lactate  
570 ligase (VanA). *Proc of the Natl Acad of Sci* **2000**. 97, 8921–8925.
- 571 [9] Case, D.; Cheatham, T.; Darden, T.; Gohlke, H.; Luo, R.; Merz, K.; Onufriev, A.;  
572 Simmerling, C.; Wang, B.; Woods, R. The Amber biomolecular simulation programs.  
573 *Journal of computational chemistry* **2005**. 26, 1668–1688.
- 574 [10] Case, D.; Cheatham, T.; Darden, T.; Gohlke, H.; Luo, R.; Merz, K.; Onufriev, A.;  
575 Simmerling, C.; Wang, B.; Woods, R. The Amber biomolecular simulation programs.  
576 *J Computat Chem* **2005**. 26, 1668–1688.
- 577 [11] Hornak, V.; Abel, R.; Okur, A.; Strockbine, B.; Roitberg, A.; Simmerling, C. Com-  
578 parison of multiple AMBER force fields and development of improved protein back-  
579 bone parameters. *Proteins* **2006**. 65, 712–725.
- 580 [12] Wang, J.; Wang, W.; Kollman, P.; Case, D. Antechamber: an accessory software  
581 package for molecular mechanical calculations. *J. Am. Chem. Soc* **2001**. 222, U403.
- 582 [13] Wang, J.; Wolf, R.; Caldwell, J.; Kollman, P.; Case, D. Development and testing of  
583 a general AMBER force field. *J Comp Chem* **2004**. 25, 1157–1174.
- 584 [14] Jorgensen, W. Quantum and statistical mechanical studies of liquids. 10. Transfer-  
585 able intermolecular potential functions for water, alcohols, and ethers. Application  
586 to liquid water. *J Am Chem Soc* **1981**. 103, 335–340.
- 587 [15] Pearlman, D.; Case, D.; Caldwell, J.; Ross, W.; Cheatham, T.; DeBolt, S.; Fer-  
588 guson, D.; Seibel, G.; Kollman, P. AMBER, a package of computer programs for  
589 applying molecular mechanics, normal mode analysis, molecular dynamics and free

- 590 energy calculations to simulate the structural and energetic properties of molecules.  
591 *Computer Physics Communications* **1995**. 91, 1–41.
- 592 [16] Berendsen, H.; Postma, J.; Gunsteren, W. V.; DiNola, A.; Haak, J. Molecular  
593 dynamics with coupling to an external bath. *J Chem Phys* **1984**. 81, 3684–3690.
- 594 [17] Nam, K.; Gao, J.; Darrin, M. An efficient linear-scaling Ewald method for long-range  
595 electrostatic interactions in combined QM/MM calculations. *J Chem Theor Comput*  
596 **2005**. 1, 2–13.
- 597 [18] Ryckaert, J.; Ciccotti, G.; Berendsen, H. Numerical integration of the cartesian  
598 equations of motion of a system with constraints: molecular dynamics of n-alkanes.  
599 *J Comp Phys* **1977**. 23, 327–341.
- 600 [19] Kohonen, T. Self-Organized formation of topologically correct feature maps **1982**.  
601 43, 59–69.
- 602 [20] Kohonen, T. *Self-Organizing Maps*. Springer Series in Information Sciences, Heidel-  
603 berg, Germany., **2001**.
- 604 [21] Kloczkowski, A.; Jernigan, R.; Wu, Z.; Song, G.; Yang, L.; Kolinski, A.; Pokarowski,  
605 P. Distance matrix-based approach to protein structure prediction. *J Struct and*  
606 *Funct Genom* **2009**. 10, 67–81.
- 607 [22] Pettersen, E.; Goddard, T.; Huang, C.; Couch, G.; Greenblatt, D.; Meng, E.; Ferrin,  
608 T. UCSF Chimera—a visualization system for exploratory research and analysis. *J*  
609 *Comput Chem* **2004**. 25, 1605–1612.
- 610 [23] <http://www.chemaxon.com/products/marvin/marvinsketch>.

- 611 [24] Shoichet, B.; Bodian, D.; Kuntz, I. Molecular docking using shape descriptors. *J*  
612 *Comput Chem* **1992**. 13, 380–397.
- 613 [25] Meng, E.; Shoichet, B.; Kuntz, I. Automated docking with grid-based energy eval-  
614 uation. *J Comput Chem* **1992**. 13, 505–524.
- 615 [26] Lang, P.; Brozell, S.; Mukherjee, S.; Pettersen, E.; Meng, E.; Thomas, V.; Rizzo, R.;  
616 Case, D.; James, T.; Kuntz, I. DOCK 6: Combining techniques to model RNA–small  
617 molecule complexes. *RNA* **2009**. 15, 1219–1230.
- 618 [27] Richards, F. Areas, volumes, packing and protein structure. *Annu Rev Biophys*  
619 *Bioeng* **1977**. 6, 151–176.
- 620 [28] Connolly, M. Solvent-accessible surfaces of proteins and nucleic acids. *Science* **1983**.  
621 221, 709–713.
- 622 [29] Kuntz, I.; Blaney, J.; Oatley, S.; Langridge, R.; Ferrin, T. A geometric approach to  
623 macromolecule-ligand interactions. *J Mol Biol.* **1982**. 161, 269–88.
- 624 [30] Srinivasan, J.; Cheatham, T.; Cieplak, P.; Kollman, P.; David, A. Continuum  
625 solvent studies of the stability of DNA, RNA, and phosphoramidate-DNA helices. *J*  
626 *Am Chem Soc* **1998**. 120, 9401–9409.
- 627 [31] Kollman, P.; Massova, I.; Reyes, C.; Kuhn, B.; Huo, S.; Chong, L.; Lee, M.; Lee, T.;  
628 Duan, Y.; Wang, W.; Donini, G.; Cieplak, P.; Srinivasan, J.; Case, D.; Cheatham, T.  
629 Calculating structures and free energies of complex molecules: combining molecular  
630 mechanics and continuum models. *Acc Chem Res* **2000**. 33, 889–897.
- 631 [32] Hawkins, G.; Cramer, C.; Truhlar, D. Pairwise solute descreening of solute charges  
632 from a dielectric medium. *Chem Phys Lett* **1995**. 246, 122–129.

- 633 [33] Hawkins, G.; Cramer, C.; Truhlar, D. Parametrized models of aqueous free energies  
634 of solvation based on pairwise descreening of solute atomic charges from a dielectric  
635 medium. *J Phys Chem* **1996**. 100, 19824–19839.
- 636 [34] Rizzo, R.; Aynechi, T.; David, A.; Kuntz, I. Estimation of absolute free energies  
637 of hydration using continuum methods: accuracy of partial charge models and opti-  
638 mization of nonpolar contributions. *J Chem Theo Comput* **2006**. 2, 128–139.
- 639 [35] Ultsch, A. Maps for the visualization of high-dimensional data spaces. In *Proc.*  
640 *Workshop on Self organizing Maps*. pages 225–230.
- 641 [36] Ultsch, A. *U\*-matrix: a tool to visualize clusters in high dimensional data*. Fach-  
642 bereich Mathematik und Informatik, **2003**.
- 643 [37] Fraccalvieri, D.; Pandini, A.; Stella, F.; Bonati, L. Conformational and functional  
644 analysis of molecular dynamics trajectories by Self-Organising Maps. *BMC bioin-*  
645 *formatics* **2011**. 12, 158.
- 646 [38] Ding, F.; Dokholyan, N.; Buldyrev, S.; Stanley, H.; Shakhnovich, E. Direct molecular  
647 dynamics observation of protein folding transition state ensemble. *Biophys J* **2002**.  
648 83, 3525–3532.
- 649 [39] Humphrey, W.; Dalke, A.; Schulten, K. VMD - Visual Molecular Dynamics. *J Mol*  
650 *Graph* **1996**. 14, 33–38.
- 651 [40] Spill, Y.; Bouvier, G.; Nilges, M. A convective replica-exchange method for sampling  
652 new energy basins. *J Comput Chem* **2013**. 34, 132–140.

- 653 [41] Bolhuis, P.; Chandler, D.; Dellago, C.; Geissler, P. Transition path sampling: Throw-  
654 ing ropes over rough mountain passes, in the dark. *Ann Rev Phys Chem* **2002**. 53,  
655 291–318.
- 656 [42] Lei, H.; Duan, Y. Improved sampling methods for molecular simulation. *Curr Opin*  
657 *Struct Biol* **2007**. 17, 187–191.
- 658 [43] Mitsutake, A.; Mori, Y.; Okamoto, Y. Enhanced sampling algorithms. *Methods Mol*  
659 *Biol* **2013**. 924, 153–195.
- 660 [44] Vanden-Eijnden, W.; Ren, W.; Vanden-Eijnden, E. Transition pathways in complex  
661 systems: Reaction coordinates, isocommittor surfaces, and transition tubes. *Chem*  
662 *Phys Lett* **2005**. 413, 242–247.
- 663 [45] Hrast, M.; Vehar, B.; Turk, S.; Konc, J.; Gobec, S.; Janezic, D. Function of the D  
664 alanine: D-alanine ligase lid loop: a Molecular Modeling and Bioactivity Study. *J*  
665 *Med Chem* **2012**. 55, 6849–6856.

## 666 List of Figures

667	1	Enzymatic reaction of a D-Ala:D-Ala ligase (Ddl), upper branch, and D-Ala:D-Lac ligase (VanA) on the lower branch. The transition state analog, phosphinate (PHY), mimics the tetrahedral intermediate $\ddagger_2$ . . . . .	33
668	2	3D X-ray crystallographic structure of VanA (PDB entry: 1E4E) colored according to its domains: the N-terminal [A2-G121] in blue, the C-terminal [G212-A342] in black, which includes the $\omega$ -loop [L236-A256] in green, and the central domain [C122-S211] in red, which includes the opposite domain [A149-Q208] in yellow. The disulfide bridge C52-C64, located in the N-terminal domain, is colored in pale blue (bottom right). . . . .	34
669	3	(a-c) Global conformational drifts, RMSD from the first conformations calculated on C $\alpha$ coordinates for: a) VanA (averaged over 8 MD trajectories), b) VanA <sub>SS</sub> (averaged over 3 MD trajectories for the red curve and over 4 trajectories for the black one), c) VanA <sub>SS</sub> .lig (averaged over 9 trajectories). (d-i) Drifts of ligases domains computed for one representative trajectory and colored according to the caption given on panel i. d) VanA, e) VanAss for the trajectory with the largest drift (VanAss_hight), f) VanA <sub>SS</sub> .lig, g) TtDdl <sub>open</sub> , h) TtDdl <sub>closed</sub> , i) TtDdl <sub>closed</sub> .lig. . . . .	35
670	4	<b>Principal Component Analysis (PCA)</b> of the C $\alpha$ dynamics covariance matrix for MD simulations run on VanA and TtDdl. The $\omega$ -loop is colored in green, its opposite domain in yellow and the remaining parts of the protein in blue. (a-e): projection of the first mode on the 3D structures, f) eigenvalues distribution. . . . .	36
671	5	<b>a)</b> U-matrix for the SOM used to analyze VanA and VanA <sub>SS</sub> trajectories. The map is toric. Labels on the U-matrix show which system mapped the different SOM areas. Black circles mark the VanA trajectory border. <b>b)</b> Projection of the RMSD values (Å) relatively to the initial conformation. Numbers show the initial conformation region (1), the main low U-matrix VanA <sub>SS</sub> basin (2) and the high U-matrix one (3). <b>c)</b> Projection of $\beta$ -strands content variation (current number minus initial one). <b>d)</b> Superposition of the first conformation (blue), of the last one (red) and of the transition states conformations (as defined by the flow analysis, yellow) of VanA <sub>SS</sub> . The colors correspond to the color-map of the RMSD projection matrix. . . . .	37
672	6	Flow analysis of the MD trajectories. VanA and VanA <sub>SS</sub> trajectories are underlined by gray circles and pink square respectively. The intersection between the two trajectories is delimited by brown diamonds. The three basins of VanA <sub>SS</sub> are numbered. The transition states ensembles of VanA <sub>SS</sub> are pointed out with black circles. Black dots stand for unvisited neurons. The color code of the arrow gives the density ( $f_{i,j}$ ) of each neuron, using the scale given at the right of the plot. The orientation of each vector indicates the resulting flow of the MD. The norms of the vectors are linked to the polarity of the corresponding flow. Zero-normed vectors are depicted by small black dots. . . . .	38
673	7	Projection of the MM-PBSA (a) and MM-GBSA (b) energies (kcal.mol <sup>-1</sup> ) on the U-matrix obtained from the molecular dynamics trajectories VanA and VanA <sub>SS</sub> . Black circles mark the VanA trajectory border. . . . .	39
674			
675			
676			
677			
678			
679			
680			
681			
682			
683			
684			
685			
686			
687			
688			
689			
690			
691			
692			
693			
694			
695			
696			
697			
698			
699			
700			
701			
702			
703			
704			
705			
706			
707			
708			
709			
710			
711			

712	8	Binding sites calculated with the SOM 3D algorithm on the run docking poses. The ligand coordinates associated with each neurone is drawn as cpk, ADP binding site is colored in cyan, D-Ala in magenta and D-Lac in purple. The figure was prepared with VMD. <sup>39</sup> . . . . .	40
713			
714			
715			
716	9	Docking of key ligands involved in the ligase mechanism. VanA and VanA <sub>SS</sub> conformations were extracted from the SOM clustering. Black circles border VanA trajectory and the areas are labeled on ATP plot. The GBSA scores expressed in kcal.mol <sup>-1</sup> were used to approximate the ligand binding free energy. . . . .	41
717			
718			
719			
720			

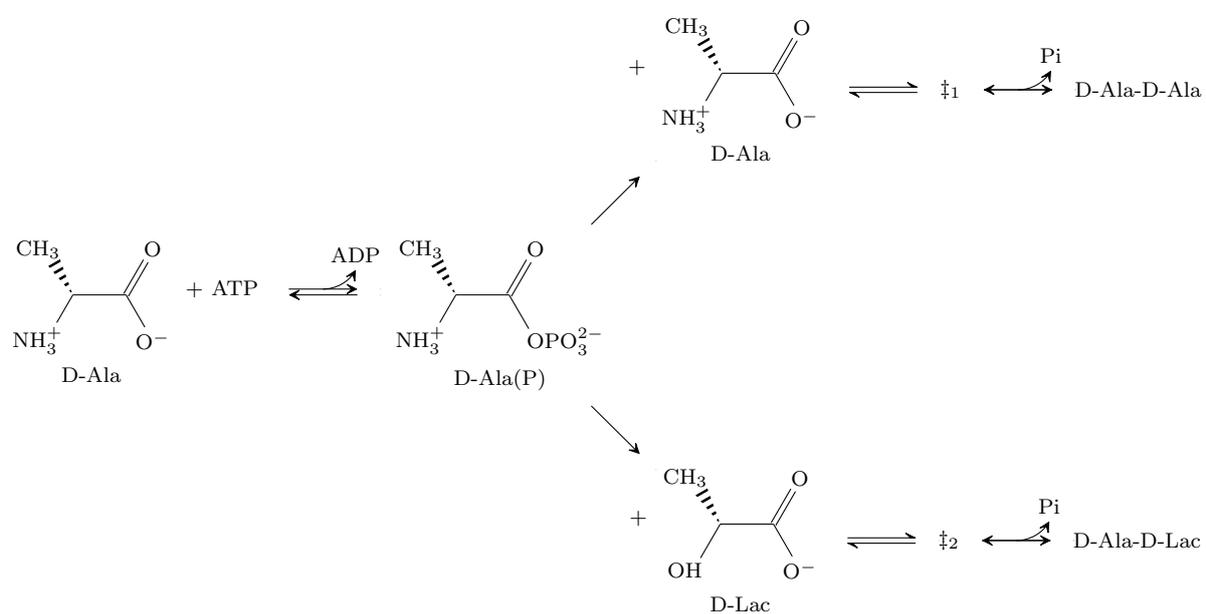


Figure 1: Enzymatic reaction of a D-Ala:D-Ala ligase (Ddl), upper branch, and D-Ala:D-Lac ligase (VanA) on the lower branch. The transition state analog, phosphinate (PHY), mimics the tetrahedral intermediate  $\ddagger_2$ .

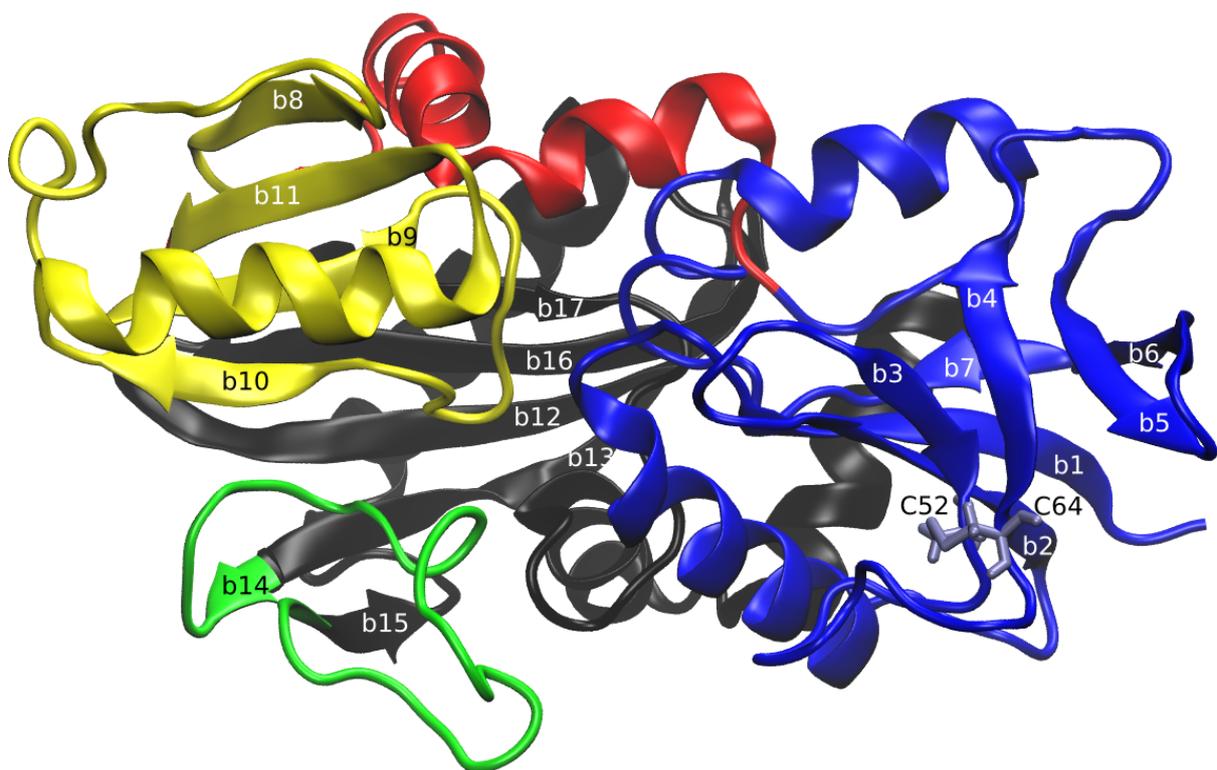


Figure 2: 3D X-ray crystallographic structure of VanA (PDB entry: 1E4E) colored according to its domains: the N-terminal [A2-G121] in blue, the C-terminal [G212-A342] in black, which includes the  $\omega$ -loop [L236-A256] in green, and the central domain [C122-S211] in red, which includes the opposite domain [A149-Q208] in yellow. The disulfide bridge C52-C64, located in the N-terminal domain, is colored in pale blue (bottom right).

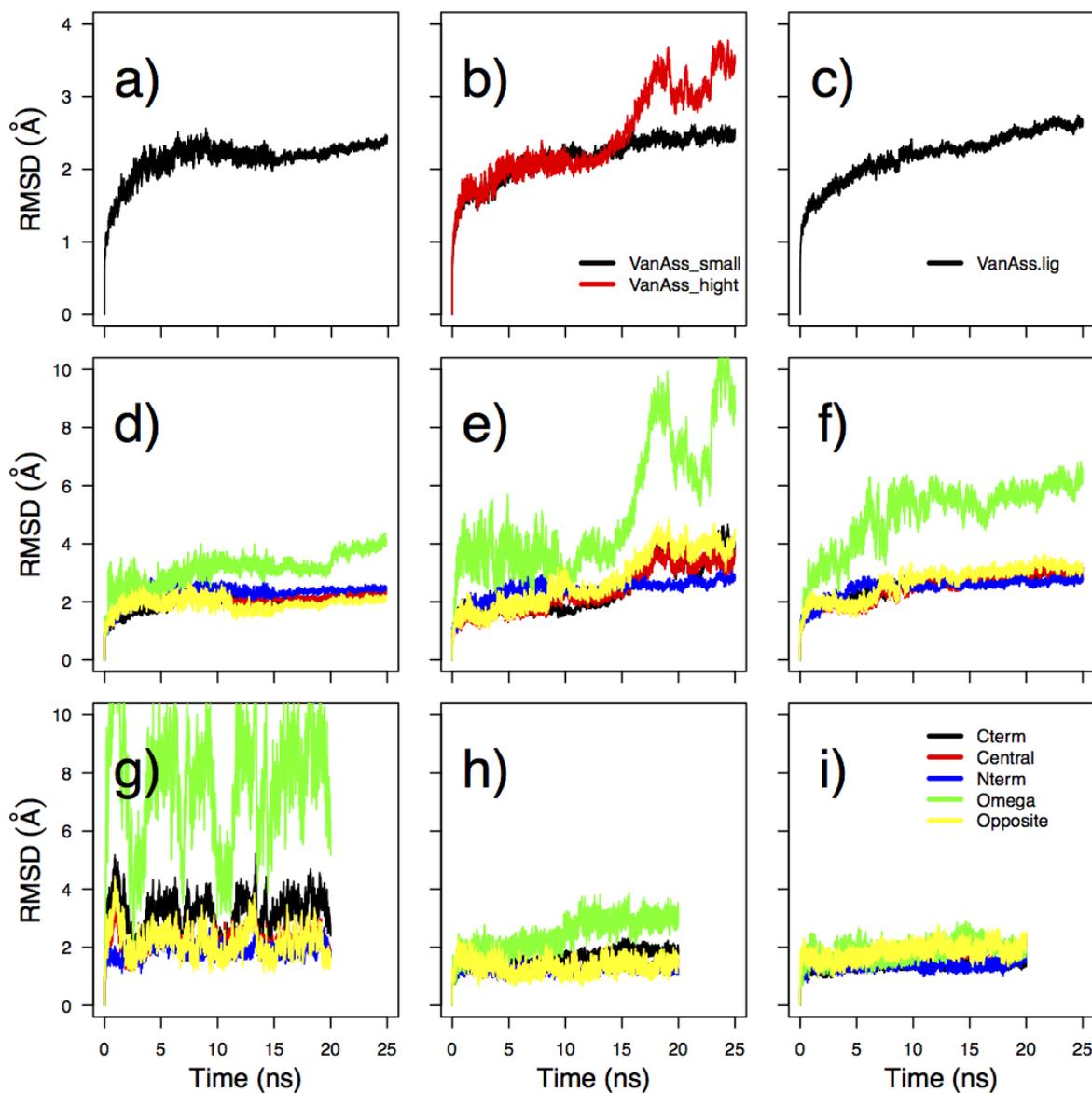


Figure 3: (a-c) Global conformational drifts, RMSD from the first conformations calculated on  $C\alpha$  coordinates for: a) VanA (averaged over 8 MD trajectories), b) VanAss (averaged over 3 MD trajectories for the red curve and over 4 trajectories for the black one), c) VanAss.lig (averaged over 9 trajectories). (d-i) Drifts of ligases domains computed for one representative trajectory and colored according to the caption given on panel i. d) VanA, e) VanAss for the trajectory with the largest drift (VanAss\_high), f) VanAss.lig, g) TtDdl<sub>open</sub>, h) TtDdl<sub>closed</sub>, i) TtDdl<sub>closed.lig</sub>.

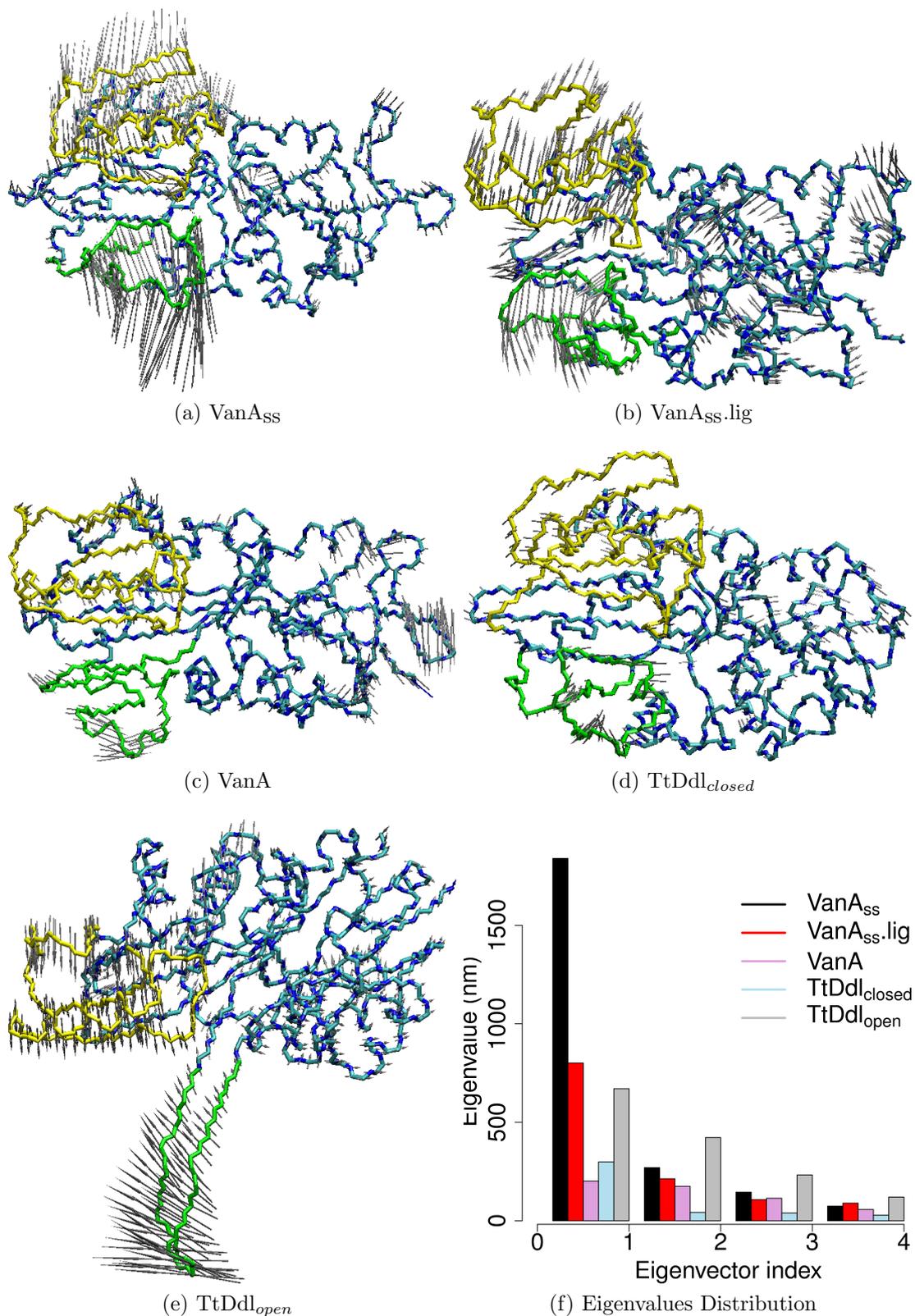


Figure 4: **Principal Component Analysis (PCA)** of the C $\alpha$  dynamics covariance matrix for MD simulations run on VanA and TtDdl. The  $\omega$ -loop is colored in green, its opposite domain in yellow and the remaining parts of the protein in blue. (a-e): projection of the first mode on the 3D structures, f) eigenvalues distribution.

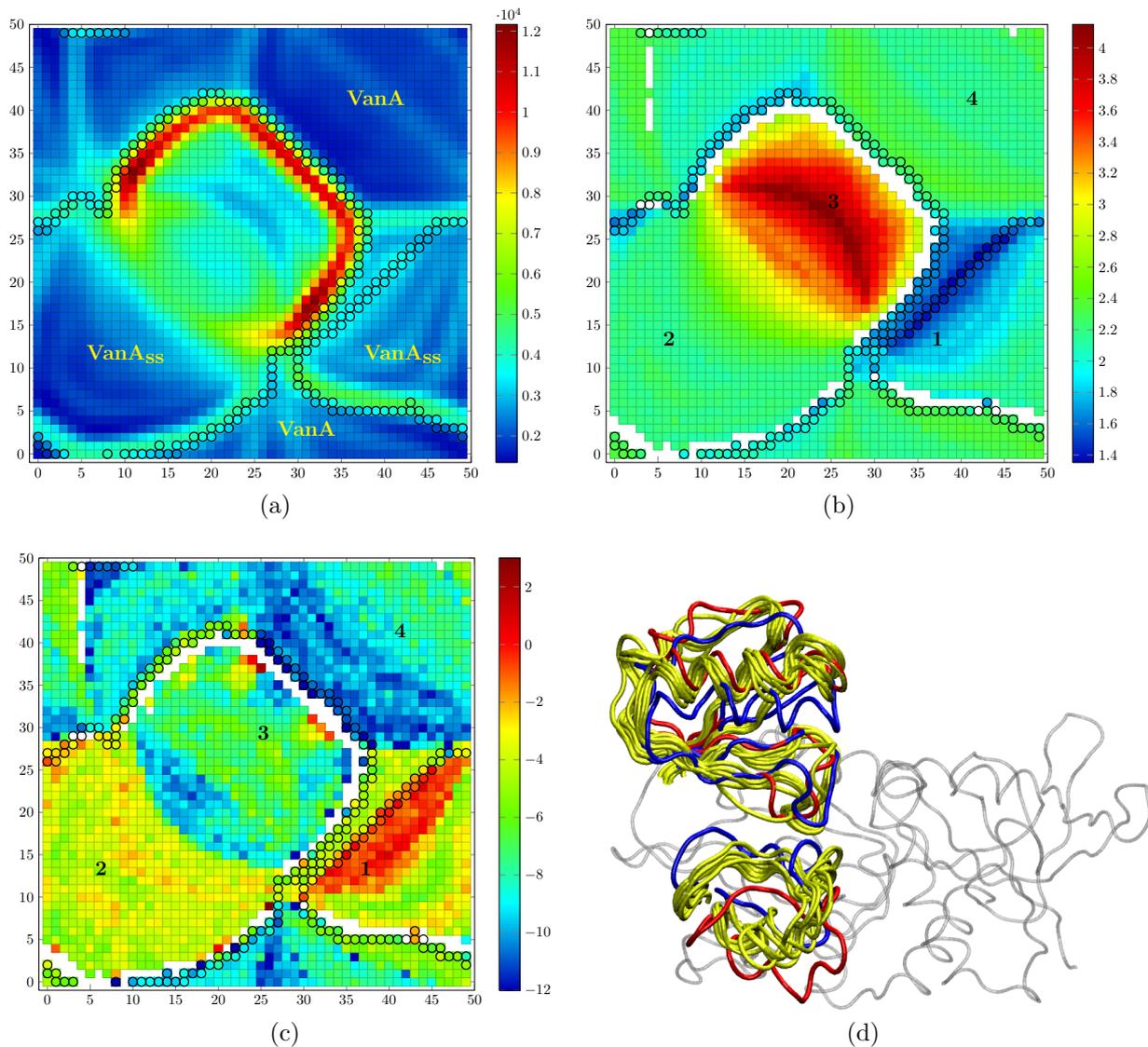


Figure 5: **a)** U-matrix for the SOM used to analyze VanA and VanA<sub>SS</sub> trajectories. The map is toric. Labels on the U-matrix show which system mapped the different SOM areas. Black circles mark the VanA trajectory border. **b)** Projection of the RMSD values (Å) relatively to the initial conformation. Numbers show the initial conformation region (1), the main low U-matrix VanA<sub>SS</sub> basin (2) and the high U-matrix one (3). **c)** Projection of  $\beta$ -strands content variation (current number minus initial one). **d)** Superposition of the first conformation (blue), of the last one (red) and of the transition states conformations (as defined by the flow analysis, yellow) of VanA<sub>SS</sub>. The colors correspond to the color-map of the RMSD projection matrix.

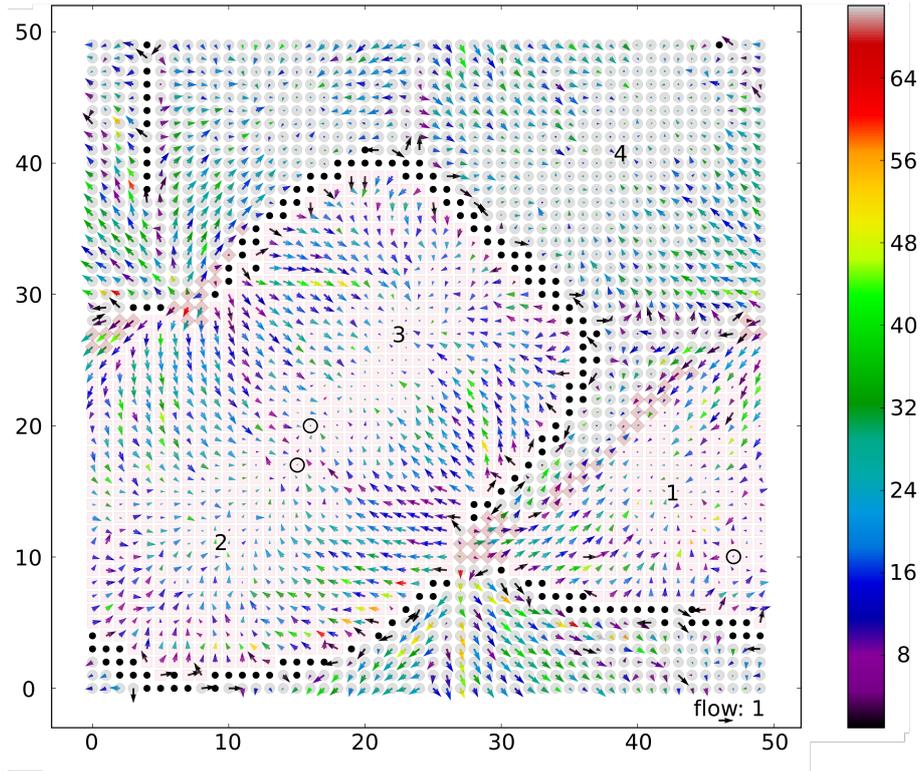


Figure 6: Flow analysis of the MD trajectories. VanA and VanA<sub>SS</sub> trajectories are underlined by gray circles and pink square respectively. The intersection between the two trajectories is delimited by brown diamonds. The three basins of VanA<sub>SS</sub> are numbered. The transition states ensembles of VanA<sub>SS</sub> are pointed out with black circles. Black dots stand for unvisited neurons. The color code of the arrow gives the density ( $f_{i,j}$ ) of each neuron, using the scale given at the right of the plot. The orientation of each vector indicates the resulting flow of the MD. The norms of the vectors are linked to the polarity of the corresponding flow. Zero-normed vectors are depicted by small black dots.

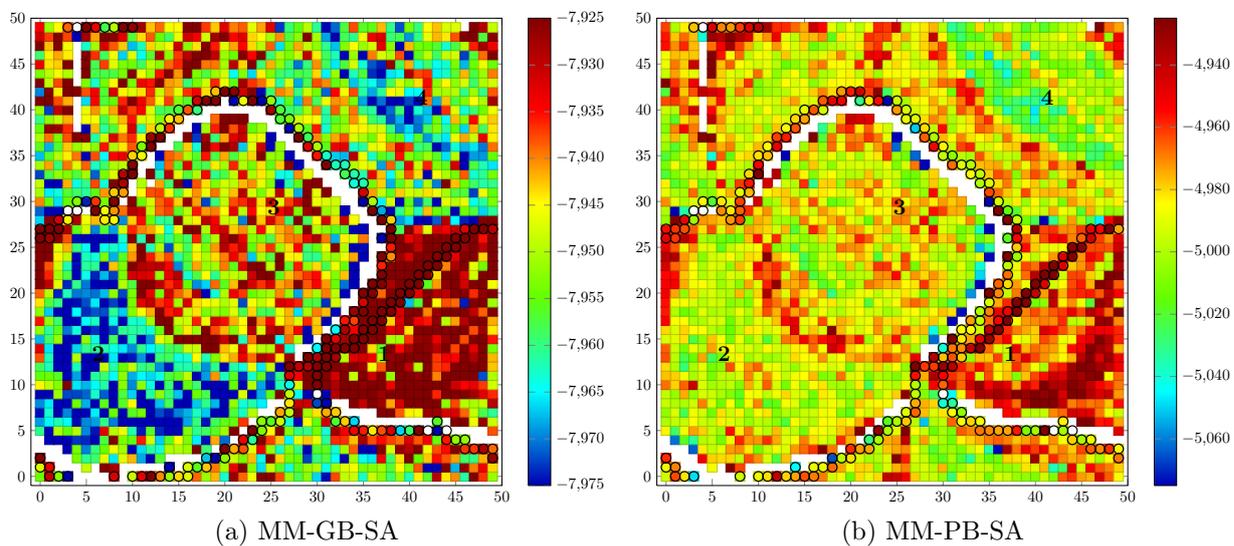


Figure 7: Projection of the MM-PBSA (a) and MM-GBSA (b) energies ( $\text{kcal.mol}^{-1}$ ) on the U-matrix obtained from the molecular dynamics trajectories VanA and VanA<sub>SS</sub>. Black circles mark the VanA trajectory border.

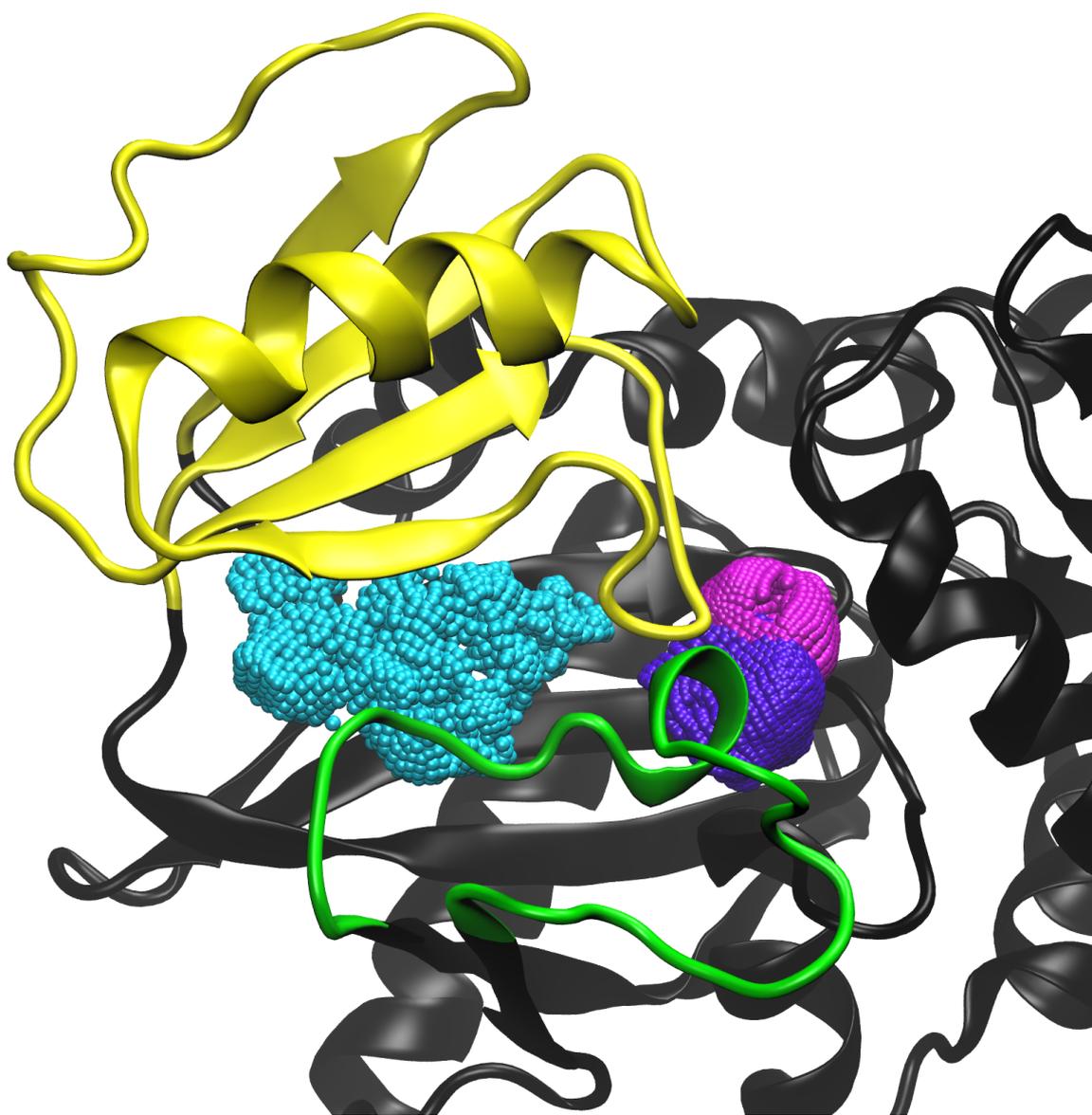


Figure 8: Binding sites calculated with the SOM 3D algorithm on the run docking poses. The ligand coordinates associated with each neurone is drawn as cpk, ADP binding site is colored in cyan, D-Ala in magenta and D-Lac in purple. The figure was prepared with VMD.<sup>39</sup>

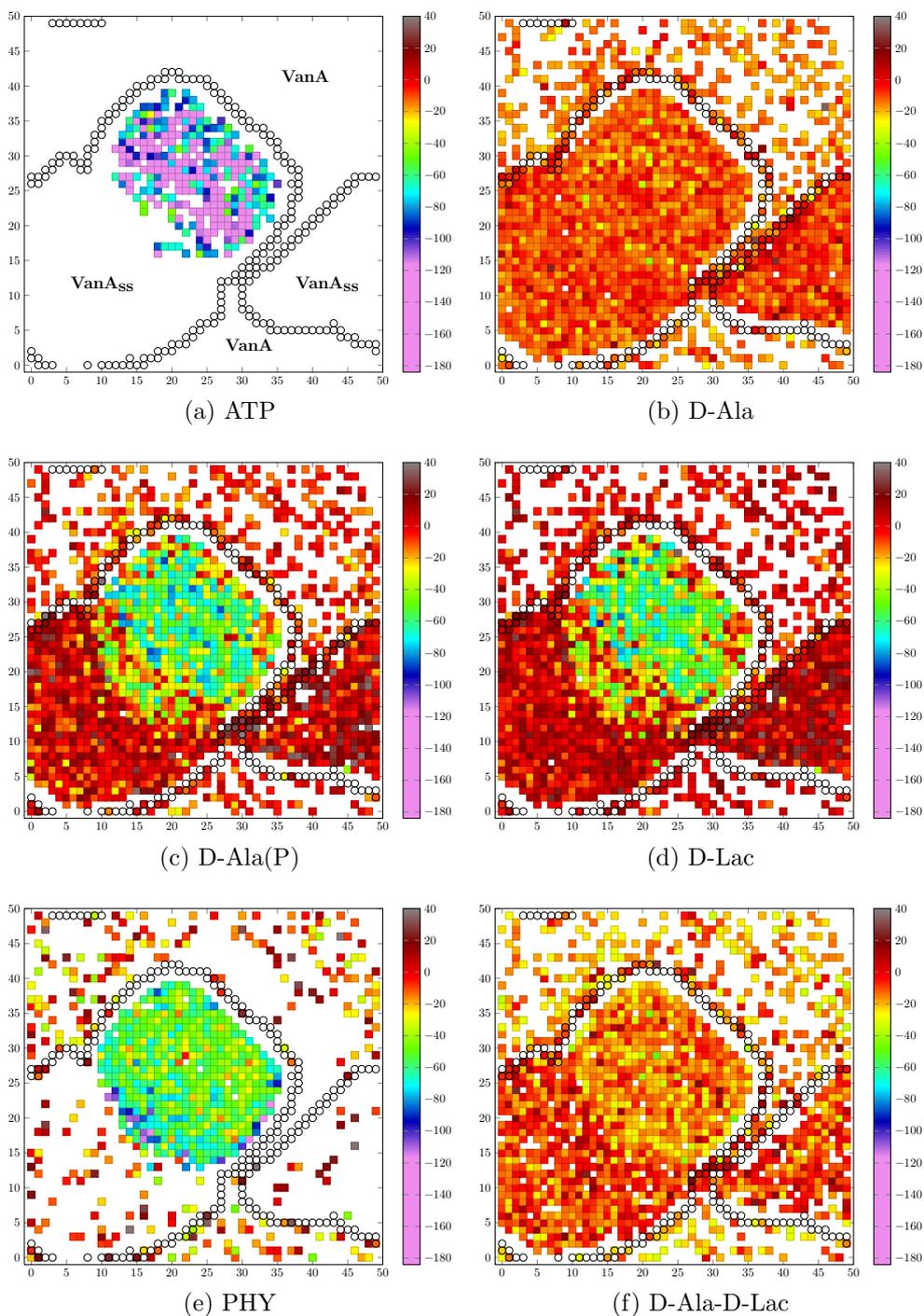


Figure 9: Docking of key ligands involved in the ligase mechanism. VanA and VanA<sub>SS</sub> conformations were extracted from the SOM clustering. Black circles border VanA trajectory and the areas are labeled on ATP plot. The GBSA scores expressed in kcal.mol<sup>-1</sup> were used to approximate the ligand binding free energy.

721 **List of Tables**

722	1	Systems used for MD simulations . . . . .	43
723	2	X-ray structures used in MD simulations . . . . .	44

PDB	Name	Ligands	Counterions	Number of water molecules	Number of recorded trajectories
2YZG	TtDdl <sub>open</sub>	-	12 Na <sup>+</sup>	13366	1
2ZDH	TtDdl <sub>closed</sub>	-	12 Na <sup>+</sup>	10854	1
2ZDH	TtDdl <sub>closed</sub> .lig	ADP, D-Ala, 2 Mg <sup>2+</sup>	11 Na <sup>+</sup>	10853	1
1E4E	VanA	-	5 Na <sup>+</sup>	13585	8
1E4E	VanA.lig	ADP, PHY, 2 Mg <sup>2+</sup>	4 Na <sup>+</sup>	13582	9
1E4E	VanA <sub>SS</sub>	-	5 Na <sup>+</sup>	13585	7
1E4E	VanA <sub>SS</sub> .lig	ADP, PHY, 2 Mg <sup>2+</sup>	4 Na <sup>+</sup>	13582	9

Table 1: Systems used for MD simulations

MD simulations	X-ray crystallographic structures	
Trajectory name	PDB	Ligands in the pocket
TtDdl <sub>open</sub>	2YZG	-
TtDdl <sub>closed</sub>	2ZDH	ADP, D-Ala, 2 Mg <sup>2+</sup>
TtDdl <sub>closed</sub> .lig	2ZDH	ADP, D-Ala, 2 Mg <sup>2+</sup>
VanA	1E4E	ADP, PHY, 2 Mg <sup>2+</sup>
VanA.lig	1E4E	ADP, PHY, 2 Mg <sup>2+</sup>
VanA <sub>SS</sub>	1E4E	ADP, PHY, 2 Mg <sup>2+</sup>
VanA <sub>SS</sub> .lig	1E4E	ADP, PHY, 2 Mg <sup>2+</sup>

Table 2: X-ray structures used in MD simulations