

AlgoSB Day V : RNA-protein docking

>: command lines on the jupyter hub

>: command lines on the laptop

]: estimated running time (4 CPUs)

Get your machine number X and student number Y and from the list [here](#)

Go to <https://avicenneX.loria.fr/user/studentY/>

Download the file [here](#). Open with `tar xzf algosb_td_ICB.tgz`

exo 1: rigid docking

In this exercise, you will assess the impact of RNA and protein flexibility on results obtained by rigid docking, using as test-case a complex between the ribosomal protein TL5 and a double-strand RNA. The structures of the protein unbound, the RNA unbound and the protein-RNA complex are all experimentally known (PDB IDs 2J01, 364D, 1FEU).

You will first run an artificial docking test, using the **bound** structures of the RNA and protein (i.e the same structure as in the experimental complex). You will then run a real-case docking, using the **unbound** structures of RNA and protein. The difference between the bound and unbound structures correspond to conformational changes that the RNA and protein undergo prior to or during the binding process.

1. Bound docking

go on: <http://www.attract.ph.tum.de/services/ATTRACT-devel/standard.html> and fill up the parameters below. The pdb files to be uploaded are in inputs/exo1/.

Partner (left menu)

- **Receptor:** structure file: upload protein_b.pdb
RMSD calculation : on
Reference RMSD PDB file: upload protein_b.pdb
- **Ligand:** structure file: upload rna_b.pdb
RMSD calculation : on
Reference RMSD PDB file: upload rna_b.pdb
What kind of molecule are you docking? : choose

Analysis: Calculate interface RMSD after docking: on
Number of structures to collect as PDB file : 20
Maximum number of structures to analyze : 1000

Computation: Name of the docking run : *bound*

Nb of CPU : **4** if you run it via jupyter, your number of CPU if you run it

on your machine

Click on **Get configuration** (bottom left). Download bound.tgz, upload it on the jupyter hub. Open a bash terminal : *New > terminal*

Run the docking by typing `$ALGOSB/scripts/run.sh bound [3']`

While the docking run, you can prepare the next docking in exo 2.

When the computation is done, you should see on the terminal a list of best solutions with their iRMSD and ranks. Download plot-result.png and out_bound-top20.pdb from server. Open plot-result.png. **How is the sampling¹? The scoring²?**

Run: `pymol out_bound-top10.pdb receptor.pdb ligand.pdb`

Check the consistency of the iRMSD values by comparing the models obtained to the bound molecules.

2. Unbound docking

Close and re-open the web-interface (not just refresh, else some parameters are not re-initialized !). Perform the docking with the unbound forms (protein/rna_ **ub**.pdb). As you still want to compare your results with the experimental structure of the complex, keep the bound form of each molecule (protein/rna_ **b**.pdb) as the Reference RMSD PDB file to compute the iRMSD. Download, run the docking and analyze the results like previously. [3']

How are the sampling and scoring affected by the molecule flexibility ?

exo 2: Flexible docking

To account for molecules flexibility during docking, we will test and compare two forms of flexible docking and their combinations:

- Pre-compute **harmonic modes** (i.e. energetically favorable directions of deformation) based on the structure of each molecule and its internal atomic forces. The modes are then used as additional degrees of freedom along the minimization. This is very demanding in term of computational time, so here we will limit ourselves to up to 4 modes, while in principle a dozen is recommended.
- Use a **conformational ensemble**, i.e. a set of different unbound structures. Those can come from different X-ray experiment, or NMR experiment, or can be obtained by MD simulations (cf lecture by S. Pasquali). Here we provide a set of 10 RNA structures coming from NMR, and 3 protein structures coming from X-ray.

We will use as a test-case a complex between the transcription factor NF-kappaB(p50) and a kinked stem-loop RNA. The structures of the unbound protein, unbound RNA and

1 Capacity of the docking to find at least one near-native structure, e.g. with low iRMSD. $\leq 1 \text{ \AA}$ is excellent, $\leq 2 \text{ \AA}$ is good, $\leq 4 \text{ \AA}$ is acceptable.

2 Capacity of the scoring function to discriminate (give low rank to) the near-native structures

protein-RNA complex are all experimentally known (PDB IDs 1LES, 2JWV, 1OOA).

On the ATTRACT web-interface, use protein_ub.pdb and rna_ub.pdb in inputs/exo2/ as receptor and ligand. Perform two docking runs:

1. Rigid docking

No flexibility, same as previously. This will be our reference test to assess the improvement obtained by using flexible docking. [3-4']

2. Flexible docking

Each of you will choose one combination of flexible options, then put a cross in the corresponding box in the on-line shared table [here](#) (to avoid that all groups use the same options). Use only one type of flexibility per molecule. Same settings as exo1 in **Analyses**.

- **“hm” : harmonic modes** (for receptor and/or ligand)
 Generate harmonic modes : on
 number of modes: 1-2 if hm on both molecules, ≤ 4 else-wise (to save time)
- **“ens” : conformational ensemble** (for receptor and/or ligand)
Receptor/Ligand: Use [protein/rna]_ub_ens[3/10].pdb

[3-7', more if hm on both molecules]

!!! Using hm on the receptor is very time consuming. If you want to use it on the protein and not the RNA, better use the RNA as receptor and the protein as ligand. Use it on both molecules only if you run the docking on your laptop on 8 or more CPUs.

While the docking runs, you can start preparing (but not run yet!) the docking in exo3.

Download and analyze the results as previously. Fill up the shared table. Remove the cross and add the data on your best-RMSD solution: iRMSD (rank). It should look like this:

	protein rigid	protein ens	protein 1 hm	protein 2 hm	protein 3 hm	protein 4 hm
RNA rigid		X			0.5 (2)	
RNA ens		1.2 (20)	X			
RNA 1 hm	X	5.2 (20)	> 10' on 4 CPU		Run on your laptop if > 4 CPU	
RNA 2 hm						
RNA 3 hm			Run on your laptop if > 4 CPU			
RNA 4 hm	X					

Does the flexible docking improve the sampling ?

exo 3: Fragments docking

In the case of single-stranded RNA, the RNA unbound form is too flexible to be observable experimentally. We cannot use classical unbound-based docking as previously. We will use a fragment-based approach to dock the RNA from its sequence. Here, we assume that we know that the RNA binds in a single-stranded state, as most protein families that bind ssRNA are well identified.

We will use as an example a poly-A RNA of 8 nucleotides binding to a 2-domains poly-A-binding protein. The structure of the bound complex has been solved experimentally (PDB ID 1CVJ). We will dock AAA fragments, compare the poses to each 3-nucleotide part of the bound RNA, then assemble chains of poses overlapping by 2 nucleotides (so 6 poses in total, to get an 8-nucl chain).

```
1           8
A A A A A A A A   bound RNA
[-----] frag1    docking poses
 [-----] frag2
      ...
          [-----] frag6
```

One should in principle use a full-size fragment library (conformational ensemble) of thousands of conformers, to cover the structural diversity of a trinucleotide. To make computation times compatible with the time constraints of this workshop, we will use a reduced ensemble of 12 conformers.

1. Fragments docking

On the web-interface, use one of the protein_ub-1.pdb as receptor and rna_ens12.pdb as ligand, in inputs/exo3/. Do **not** use any harmonic mode. Deactivate RMSD calculation. Run the docking [2-5 min].

After docking is done, first check if each fragment has been correctly sampled. From the docking directory, run: [\\$ALGOSB/scripts/rmsd.sh](#). In the terminal, you will see the 5 best solutions obtained for each fragment, with its rank by energy.

2. Fragments assembly

Assembly the docking poses into chains of up to 6 fragments (up to 8 nucl) with:

[\\$ALGOSB/scripts/assemble.sh](#)

```
nfrag      number of tri-nucleotide to assemble
cutoff     overlap cutoff (in Å). Start with a small one. LIMIT: 5
nposes    number of top-ranked poses to assemble. Recommended [100 - 10000]
meanrank  max geometric mean of pose ranks in each chain ( ≤ npose)
maxchains max number of chains in output. Recommended ≤ 100000
```

```
ex: $ALGOSB/scripts/assemble.sh 6 1 100 100 1000
      nfrag cutoff npose meanrank maxchains
```

When you have found a suitable set of parameters that give you an acceptable solutions (iRMSD < 5Å):

Download `chains-*.pdb result.pdb rna_b.pdb protein_b.pdb`

Visualize with `pymol result.pdb chains-*.pdb rna_b.pdb protein_b.pdb`

Fill up the on-line shared table [here](#).

nfrag	cutoff	nposes	meanrank	maxchains	best RMSD	rank
6	1	100	100	1000	4.2	1024