

# Molecular Docking in the Cloud: Introduction to Molecular Docking

PhD. Pablo Ricardo Arantes  
PhD. Conrado Pedebos

Porto Alegre, July 14<sup>th</sup> 2025

# Molecular Docking Limitations

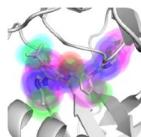
- Main limitations:
  - Rigid receptor (flexibility is possible, but limited and error-prone)
  - Computational cost for large datasets (important, but we'll talk more about this on Thursday)
  - Limited treatment of relevant effects (entropy, enthalpy, solvation)
  - **Scoring functions** (simplified approximations / many false positives and/or negatives / difficult to predict true experimental binding energies)

# AI-powered Molecular Docking

- Many new methods are being developed
- GNINA → Use of Convolutional Neural Networks (CNN) to evaluate poses
  - A fork of *smina*, which is a fork of *vina*
  - Increased accuracy for pose prediction and affinity scores
  - Trained using datasets of protein-ligand binding (PDBbind)
  - Autodock Vina is still used to generate initial poses → scoring with GNINA

## gnina/gnina

A deep learning framework for molecular docking



208  
Contributors

6  
Issues

739  
Stars

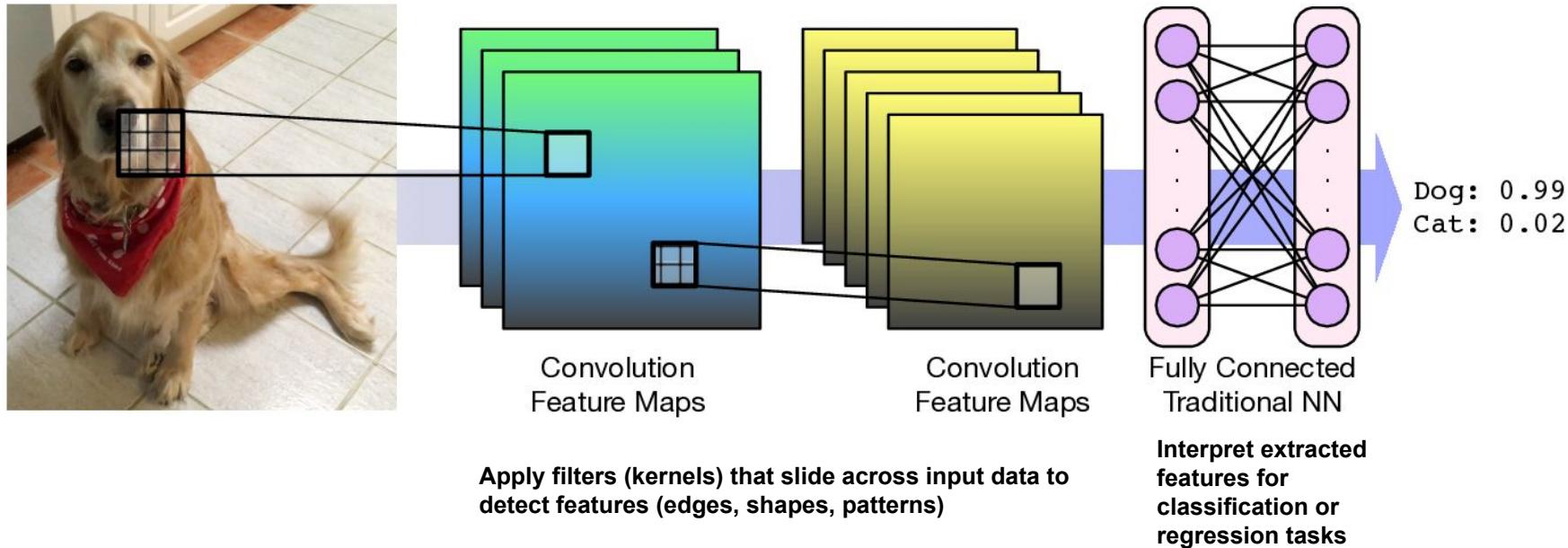
162  
Forks



<https://github.com/gnina/gnina>

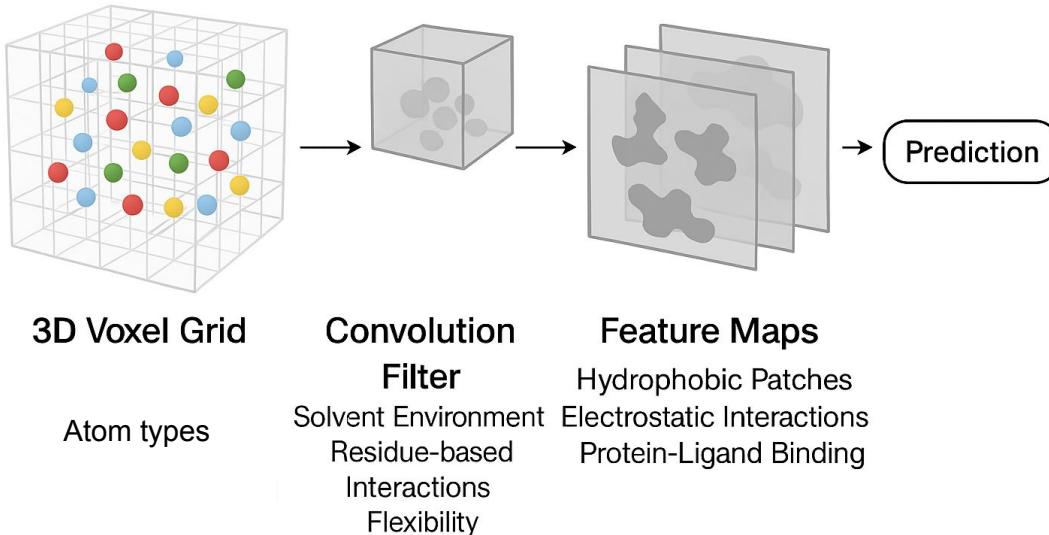
# AI-powered Molecular Docking

Convolutional neural networks learn spatially related features of an input grid to generate a prediction.



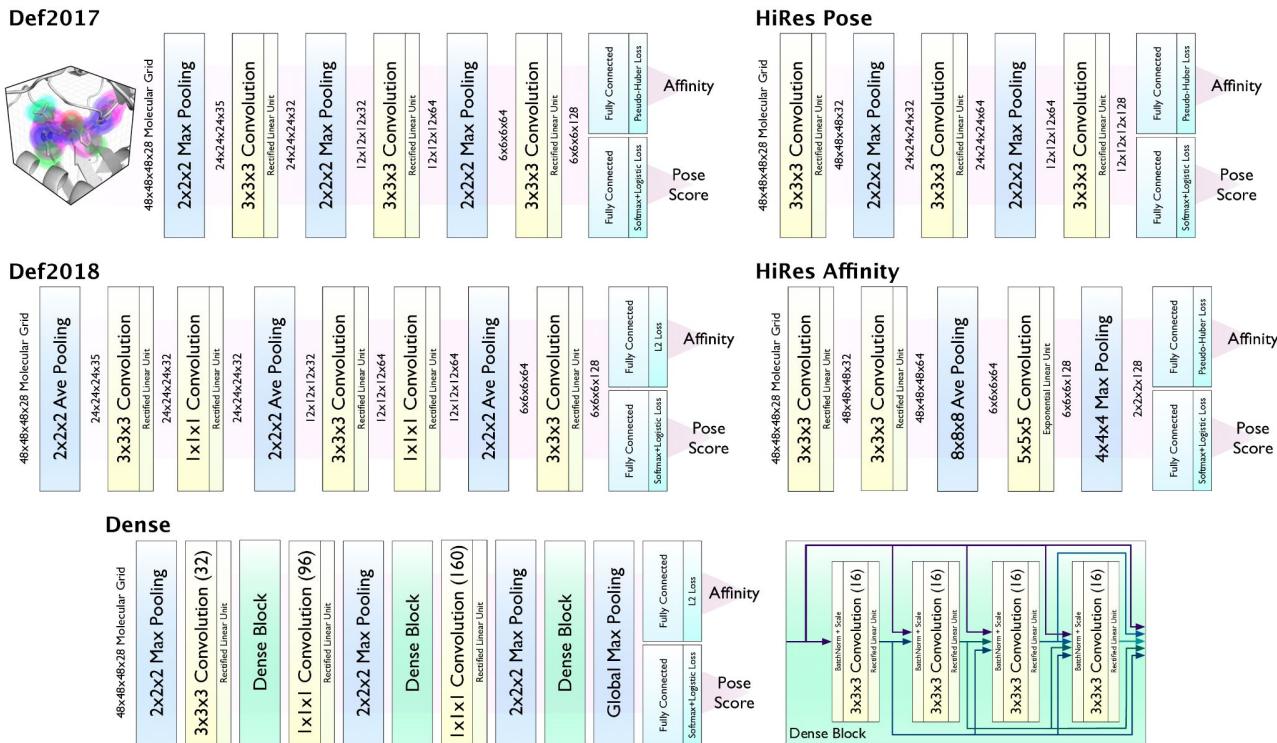
# AI-powered Molecular Docking

**Convolution neural network (CNN) in a molecular context**

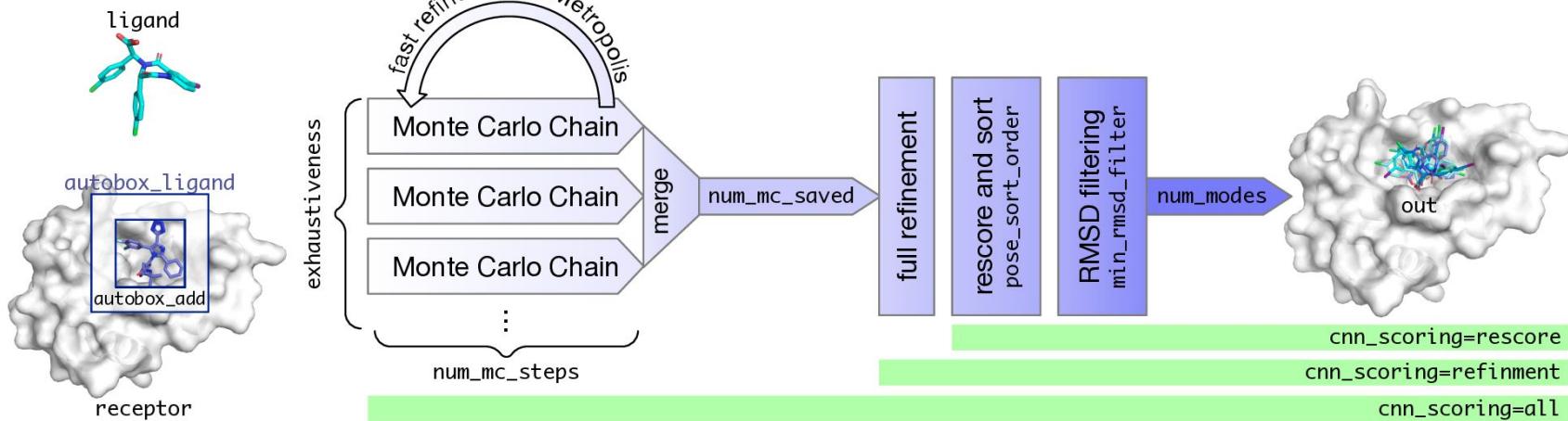


# AI-powered Molecular Docking

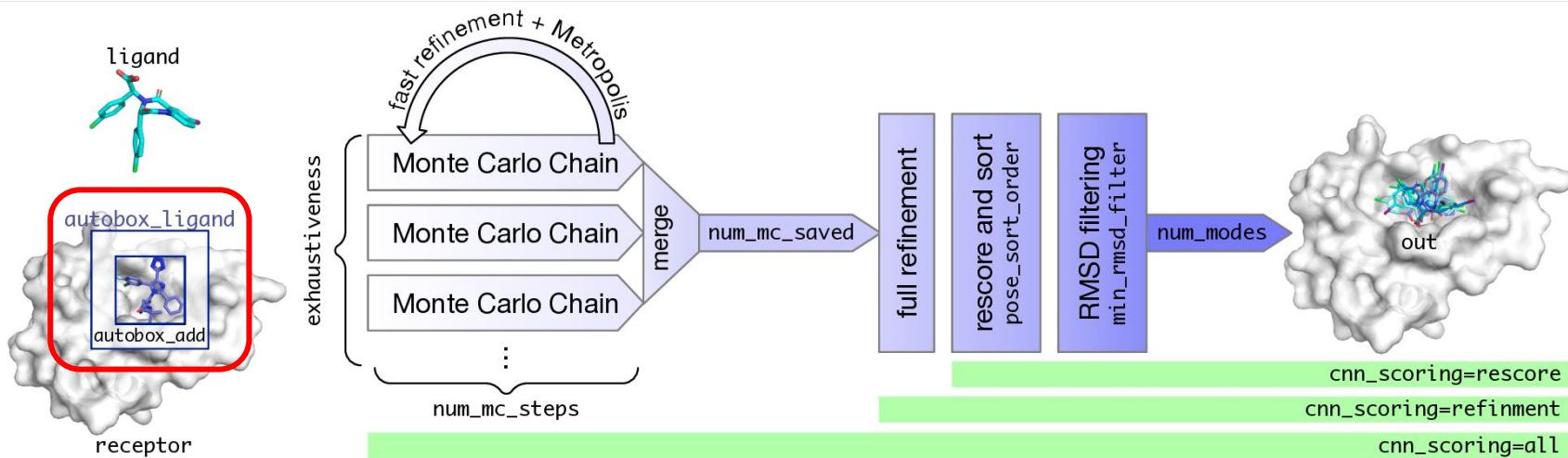
- Many different models can be employed
  - A combination (ensemble) of five models was identified to give the best performance
  - CNN model predicts both pose quality and binding affinity



# AI-powered Molecular Docking

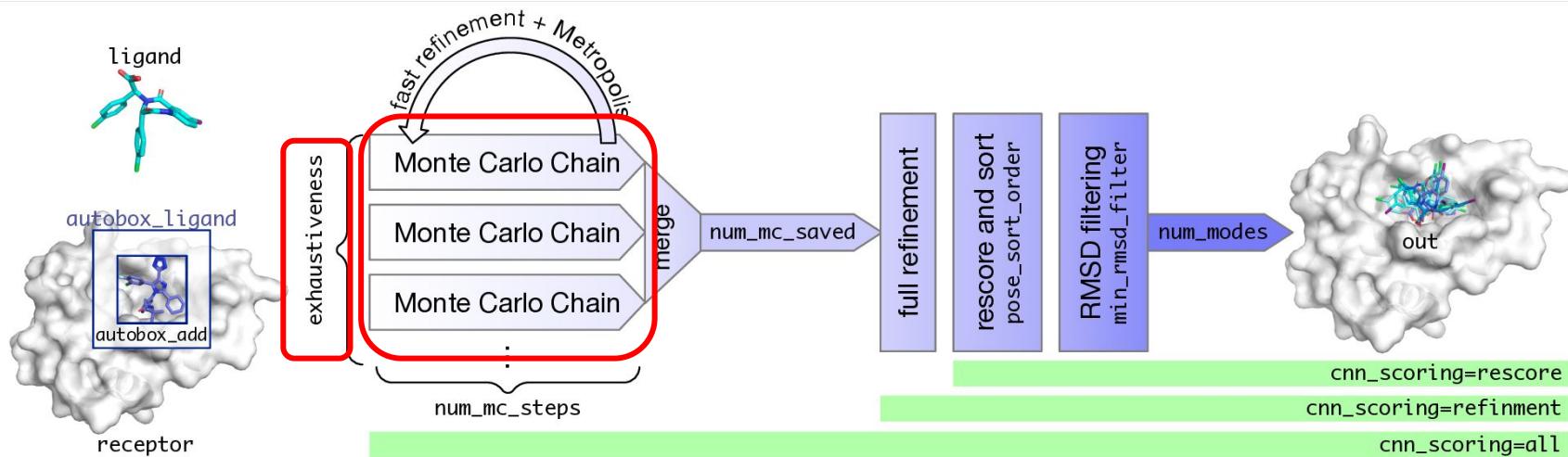


# AI-powered Molecular Docking



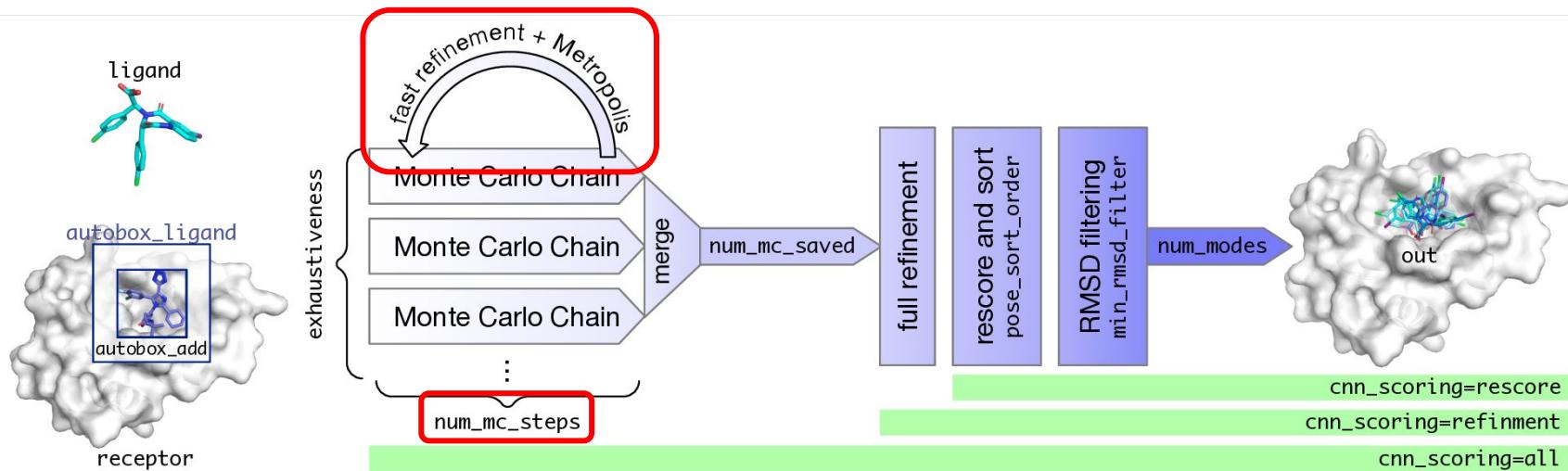
**autobox\_ligand** = creates the binding site box based on the ligand structure  
**autobox\_add** = adds the amount of buffer space

# AI-powered Molecular Docking



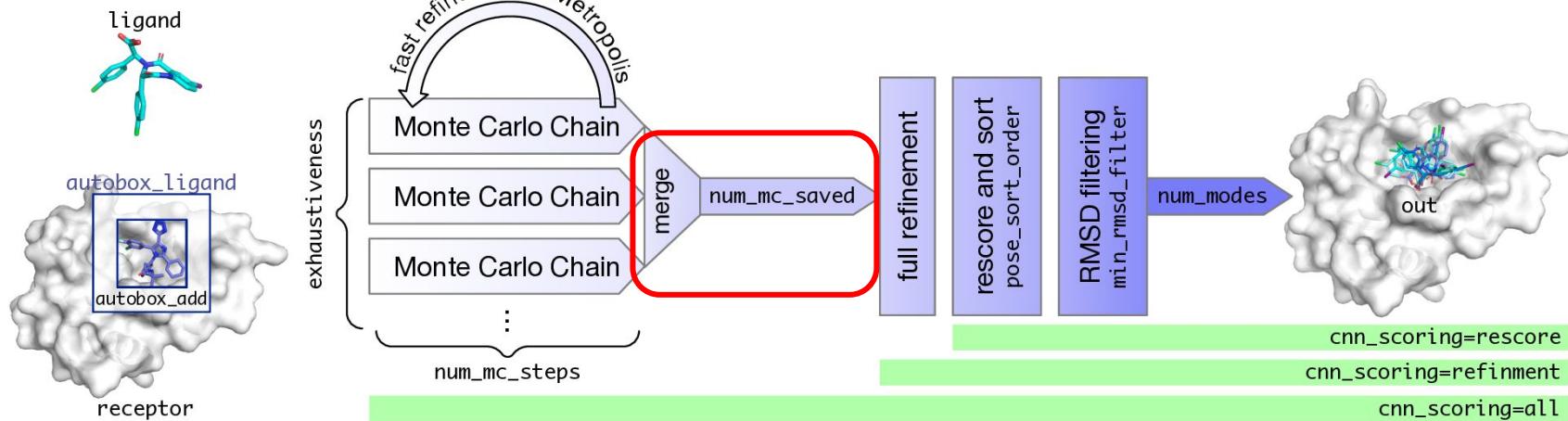
**Monte Carlo Chain = Monte Carlo sampling of the ligand conformations**  
**exhaustiveness = number of chains running concomitantly**

# AI-powered Molecular Docking



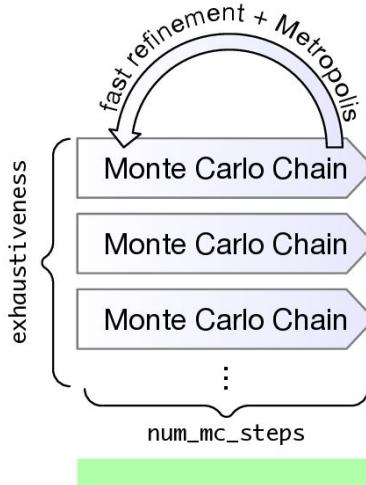
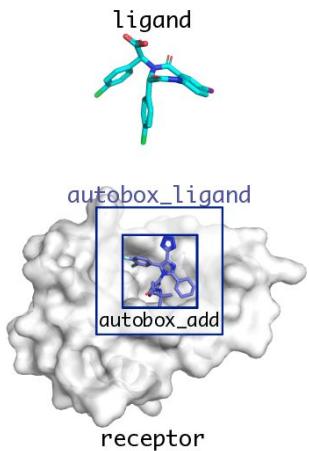
**fast refinement + Metropolis** = energy minimization and criterion for acceptance of the new state  
**num\_mc\_steps** = automatically determined (can be changed)

# AI-powered Molecular Docking

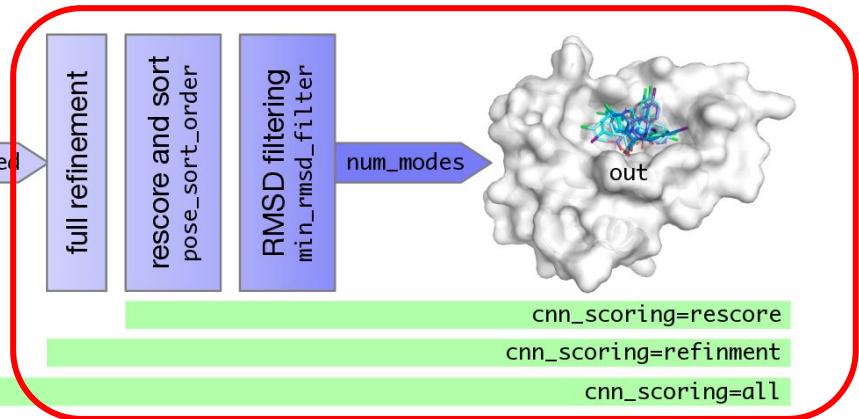


**num\_mc\_saved = 50 (default), number of top conformations retained**

# AI-powered Molecular Docking



**full refinement = shifts the ligand pose to a local energy minimum**

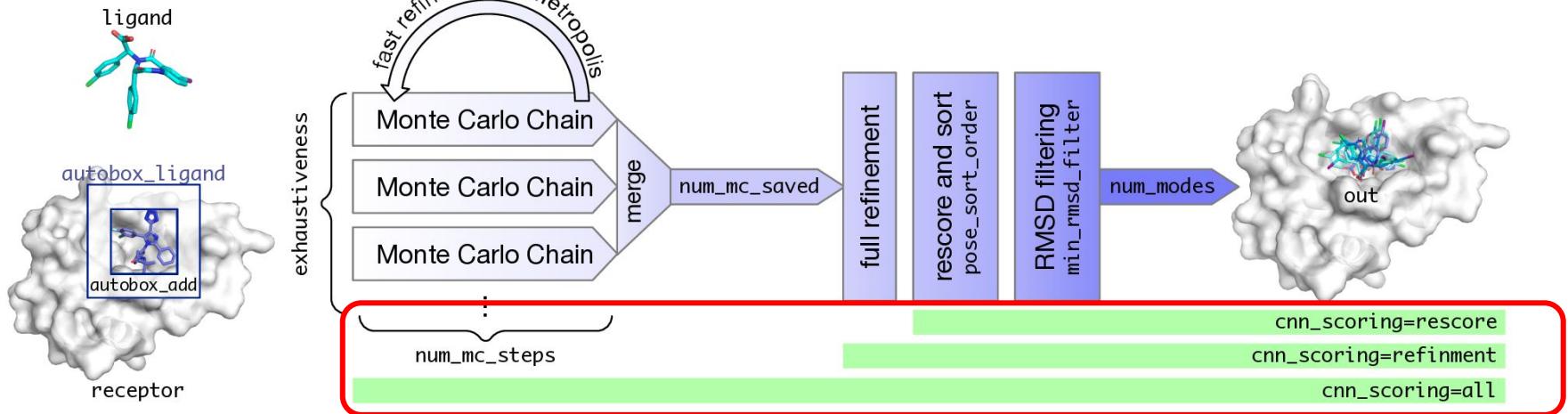


**pose\_sort\_order = how the top scoring poses are ranked**

**min\_rmsd\_filter = ligand conformations too similar are filtered out**

**num\_modes = final number of poses for analysis (default = 10)**

# AI-powered Molecular Docking



# AI-powered Molecular Docking

```
gnina --receptor 1BCU_PROT.pdb --ligand 1BCU_LIG.sdf
--out 1BCU_gnina_poses.sdf.gz --autobox_ligand
1BCU_LIG.sdf --autobox_add 4 --cnn
crossdock_default2018 dense_3 --cnn_scoring rescore
--exhaustiveness 8 --num_mc_saved 50 --cnn_rotation 0
--num_modes 9 --min_rmsd_filter 1
```

# AI-powered Molecular Docking

```
gnina --receptor 1BCU_PROT.pdb --ligand 1BCU_LIG.sdf
--out 1BCU_gnina_poses.sdf.gz --autobox_ligand
1BCU_LIG.sdf --autobox_add 4 --cnn
crossdock_default2018 dense_3 --cnn_scoring rescore
--exhaustiveness 8 --num_mc_saved 50 --cnn_rotation 0
--num_modes 9 --min_rmsd_filter 1
```

Different rotations of the protein-ligand complex the CNN is able to see for each conformation (does not seem to alter the results, so it is set to 0)

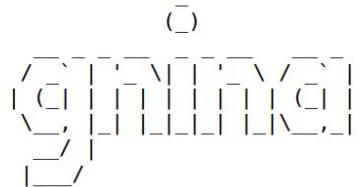
# AI-powered Molecular Docking

- What's behind the workflow?
  - OpenBabel for input parsing (PDB, sdf, mol2, etc)
  - Binding site: cartesian coordinates or automatic box around the ligand (input ligand)
    - **First lecture:** specified manually
    - **autobox\_ligand** is used to define the binding site, a rectangular prism is constructed using the minimum and maximum values for the x, y, and z coordinates of the ligand to which additional spacing (**autobox\_add**) is added in every dimension (**ligand needs to rotate freely**)
  - Scoring functions:
    - If CNN for rescoring only: user-input or built-in (Vina, Vinardo)
    - CNN function: user-input or built-in (*crossdock\_default2018, dense, general\_default2018, redock\_default2018, and default2017*) - each of which is trained using different training data and/or a different model architecture

# AI-powered Molecular Docking

- What's behind the workflow?
  - CNN calculations are performed using the cuDNN accelerated Caffe deep learning framework (GNINA 1.0) or PyTorch (GNINA 1.3)
  - CNN models predict:
    - Pose score (**CNNscore**) = probability that the pose has a low root mean square deviation (RMSD) to the binding pose
    - Binding affinity (pK) (**CNNAffinity**)
  - *cnn\_scoring default* = CNN scoring for the final resorting of ligand conformations and the empirical scoring function everywhere else in the pipeline.
  - *cnn\_scoring*:
    - “none” = smina pipeline
    - “rescoring” = CNN used in the final sorting of ligand conformations
    - “refinement” = same as the previous step but CNN is also used after Monte Carlo chain step for refining
    - “all” = CNN used at every step including energy minimization

# AI-powered Molecular Docking



```
gnina v1.0.1 HEAD:aa41230 Built Mar 23 2021.
gnina is based on smina and AutoDock Vina.
Please cite appropriately.
```

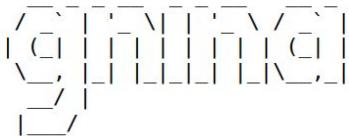
```
Commandline: ./gnina -r rec.pdb -l lig.pdb --autobox_ligand lig.pdb
Using random seed: -216854720
```

	0%	10%	20%	30%	40%	50%	60%	70%	80%	90%	100%
	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----
	*****										

mode	affinity (kcal/mol)	CNN pose score	CNN affinity
1	-8.51	0.8985	6.783
2	-8.30	0.4491	6.450
3	-6.80	0.3258	6.043
4	-7.34	0.3023	6.230
5	-5.90	0.1754	5.397
6	-6.33	0.1679	5.559
7	-6.98	0.1668	5.825
8	-5.24	0.1607	5.505
9	-7.00	0.1523	5.957

# AI-powered Molecular Docking

(



```
gnina v1.0.1 HEAD:aa41230 Built Mar 23 2021.
gnina is based on smina and AutoDock Vina.
Please cite appropriately.
```

```
Commandline: ./gnina -r rec.pdb -l lig.pdb --autobox_ligand lig.pdb
Using random seed: -216854720
```

```
0%   10    20    30    40    50    60    70    80    90   100%
|-----|-----|-----|-----|-----|-----|-----|-----|-----|
*****
```

mode	affinity (kcal/mol)	CNN pose score	CNN affinity
1	-8.51	0.8985	6.783
2	-8.30	0.4491	6.450
3	-6.80	0.3258	6.043
4	-7.34	0.3023	6.230
5	-5.90	0.1754	5.397
6	-6.33	0.1679	5.559
7	-6.98	0.1668	5.825
8	-5.24	0.1607	5.505
9	-7.00	0.1523	5.957

Affinity is the original scoring by *vina* (or *vinardo*) scoring function

CNN pose score is a value between 0 and 1 that is used to rank the poses of the ligand, where a score of 1 denotes a perfect ligand pose

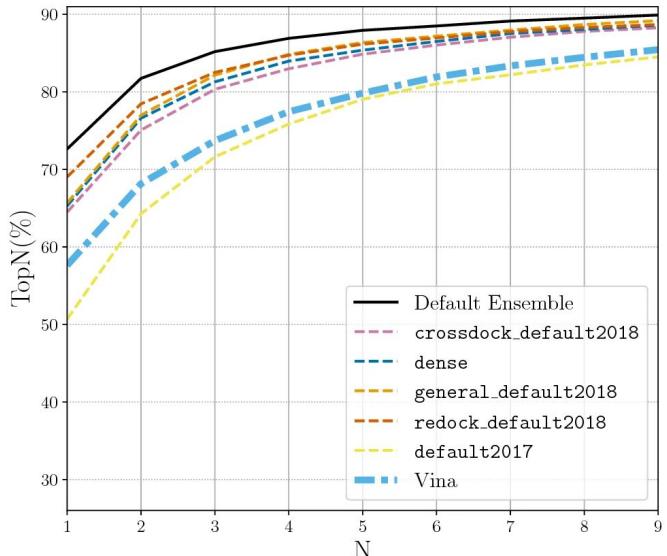
CNNaffinity is the affinity of the docked complex as determined by the CNN

# AI-powered Molecular Docking

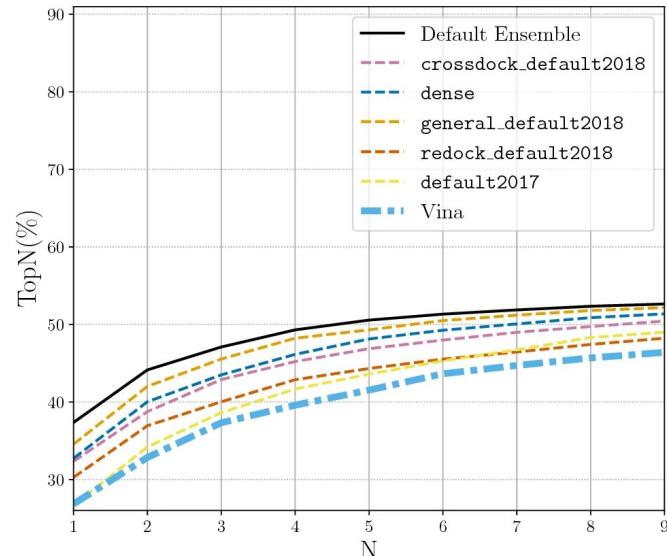
Default ensemble = *dense*,  
*general\_default2018\_3*, *dense\_3*,  
*crossdock\_default2018*, and  
*redock\_default2018*

(only “rescoring” option used)

Top1 = >70%



(a) Redocking

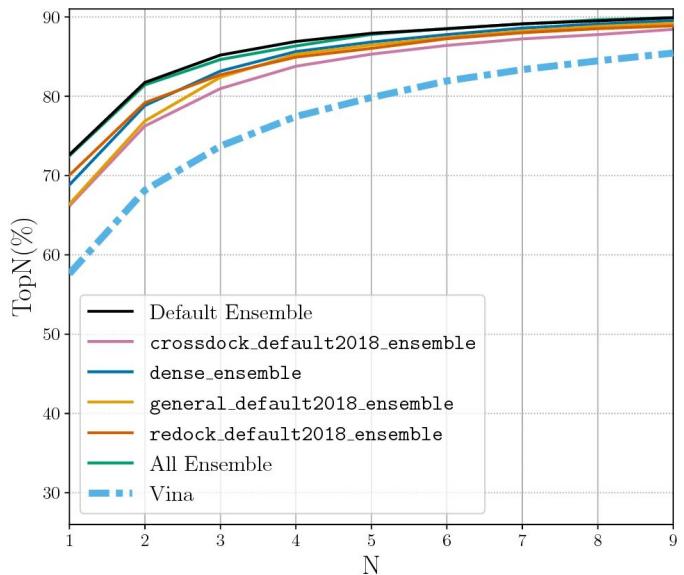


(b) Cross-docking

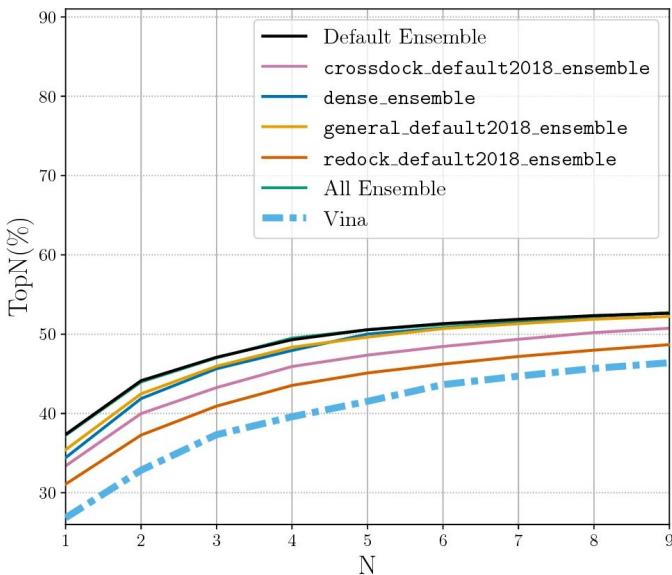
In 70% of the protein-ligand docking targets tested, the top 1 pose had an RMSD < 2.0 Å from the crystal pose.

# AI-powered Molecular Docking

Using All Ensembles achieve similar results as the Default (not worth it!)



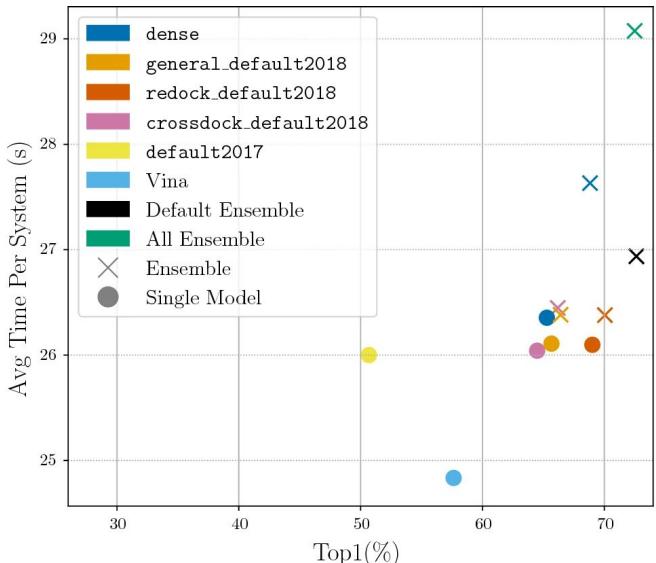
(a) Redocking



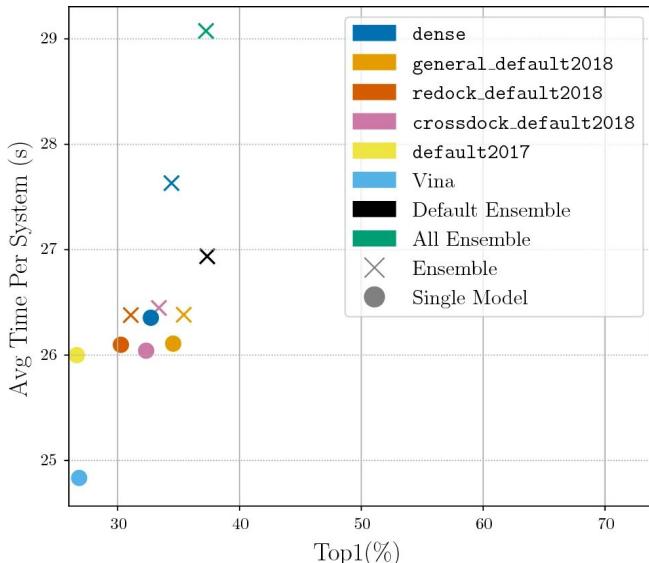
(b) Cross-docking

# AI-powered Molecular Docking

Default is not the fastest, but achieves the best results



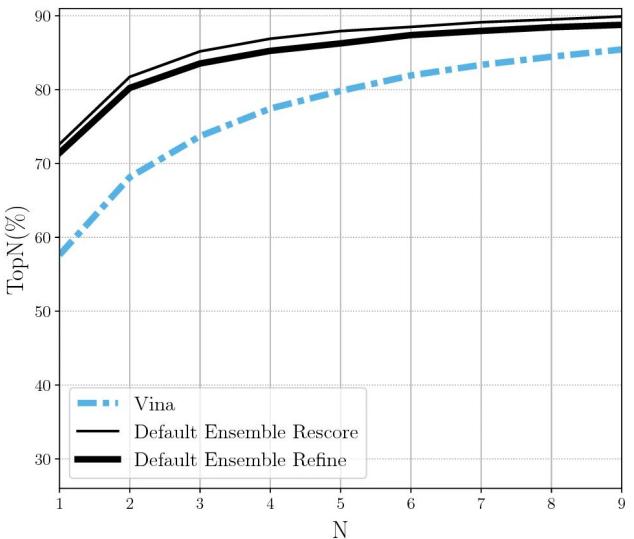
(a) Redocking



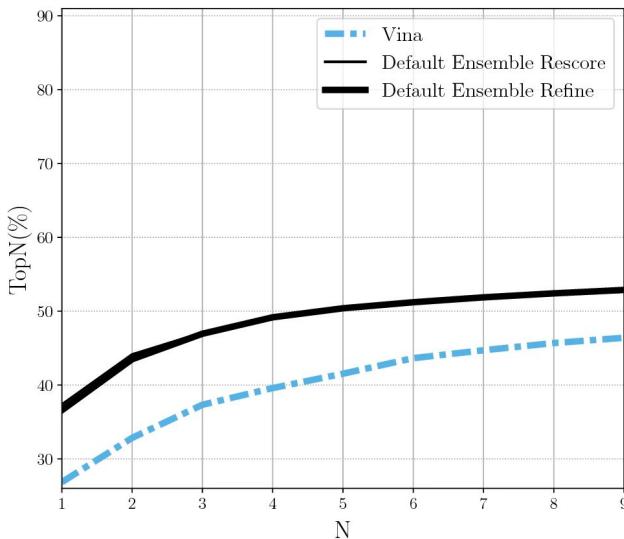
(b) Cross-docking

# AI-powered Molecular Docking

Increasing the CNN level to Refine does not achieve better results



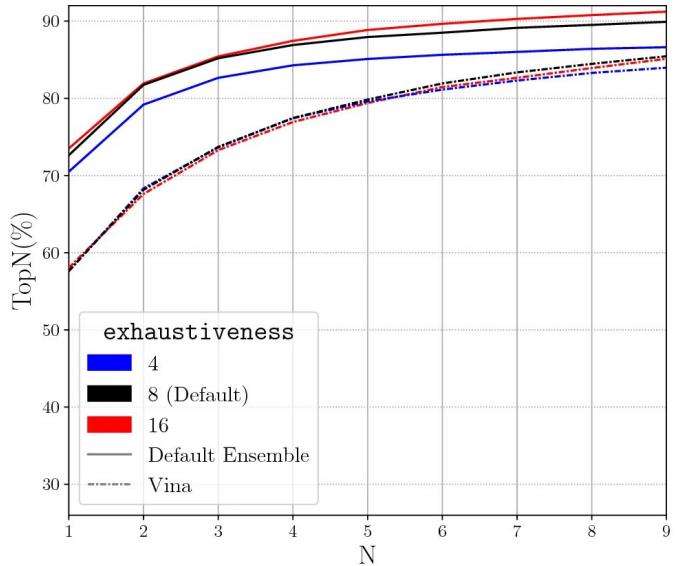
(a) Redocking



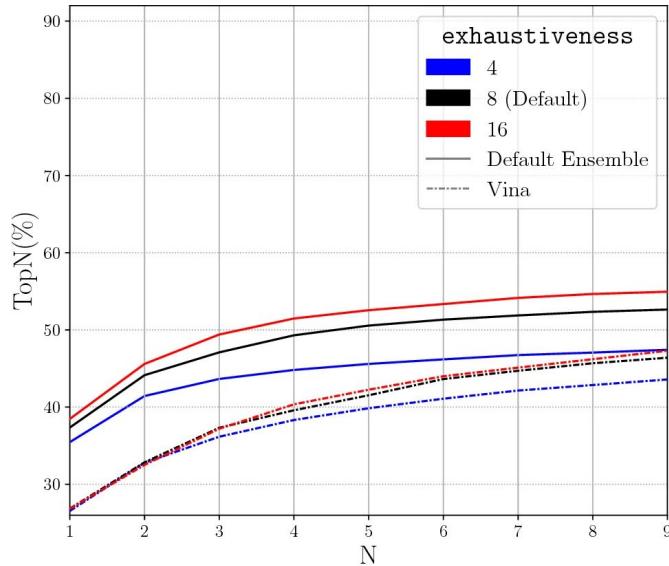
(b) Cross-docking

# AI-powered Molecular Docking

In typical docking calculations, increasing the exhaustiveness parameters yields better results, but it might not be worth it...



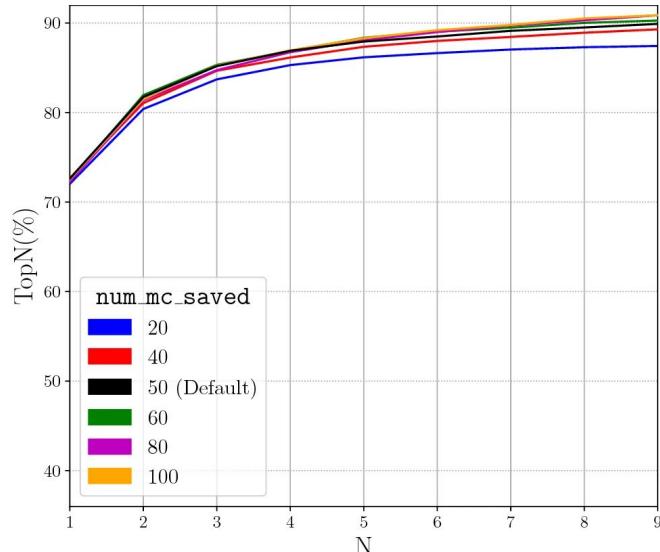
(a) Redocking



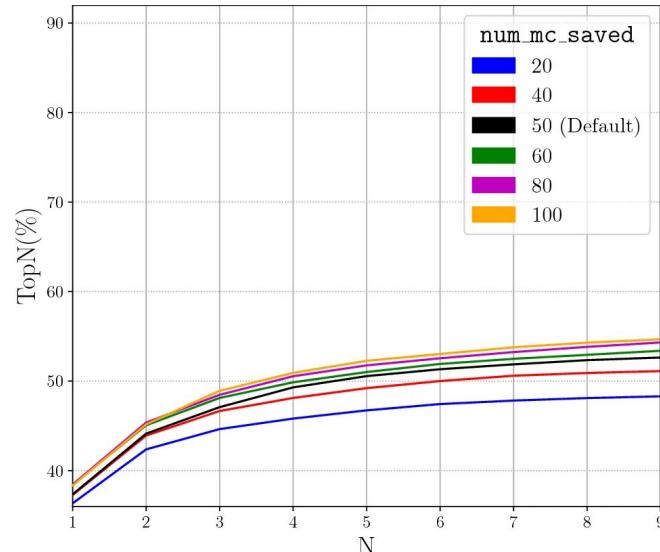
(b) Cross-docking

# AI-powered Molecular Docking

default (50) num\_mc\_saved  
is probably the best  
cost-benefit scenario



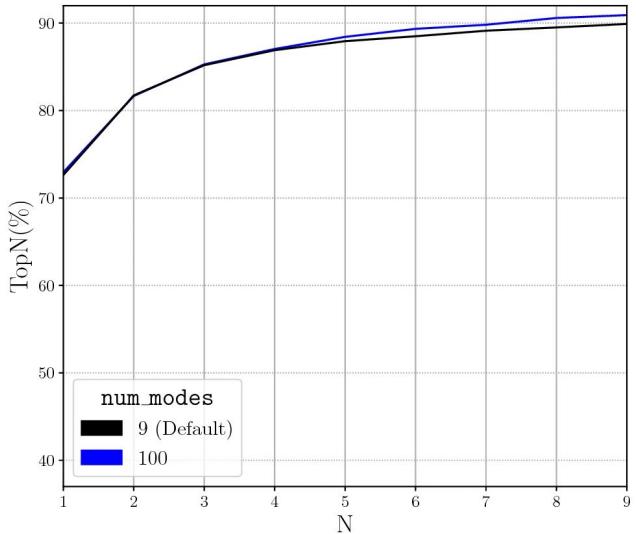
(a) Redocking



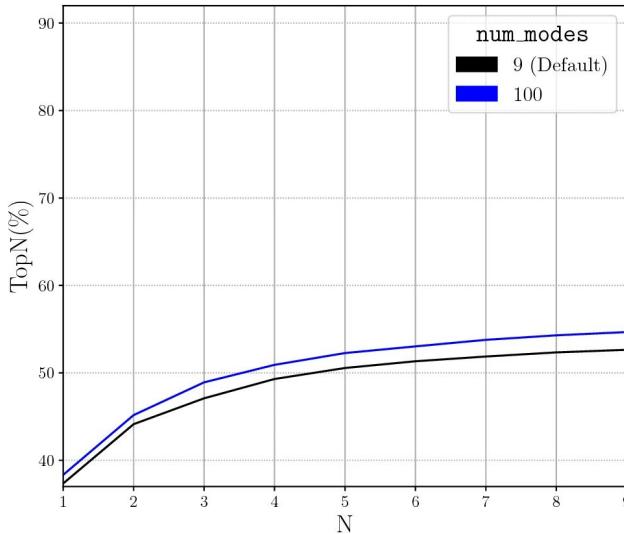
(b) Cross-docking

# AI-powered Molecular Docking

The first 9 poses are usually enough to get achieve 90% TopN (at least one pose below 2 Å of RMSD)



(a) Redocking

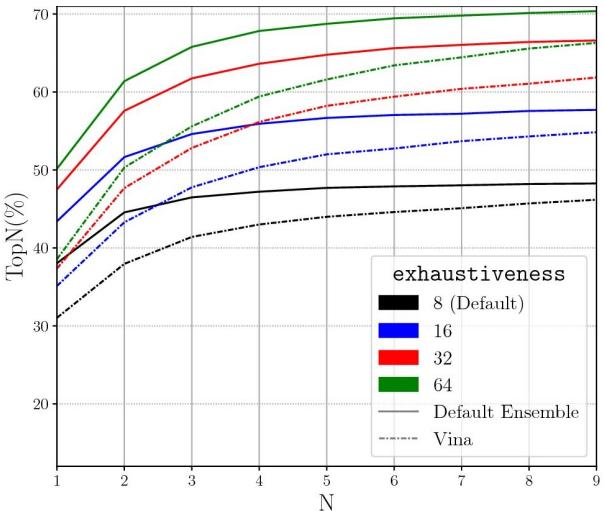


(b) Cross-docking

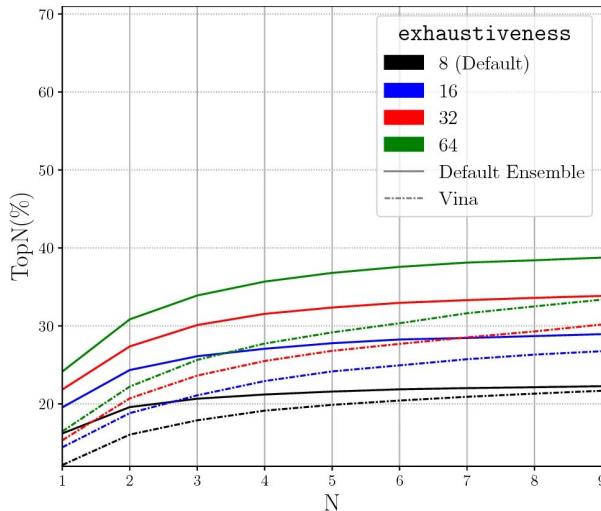
# AI-powered Molecular Docking

Things change when we do a blind docking run!

Use max. exhaustiveness available

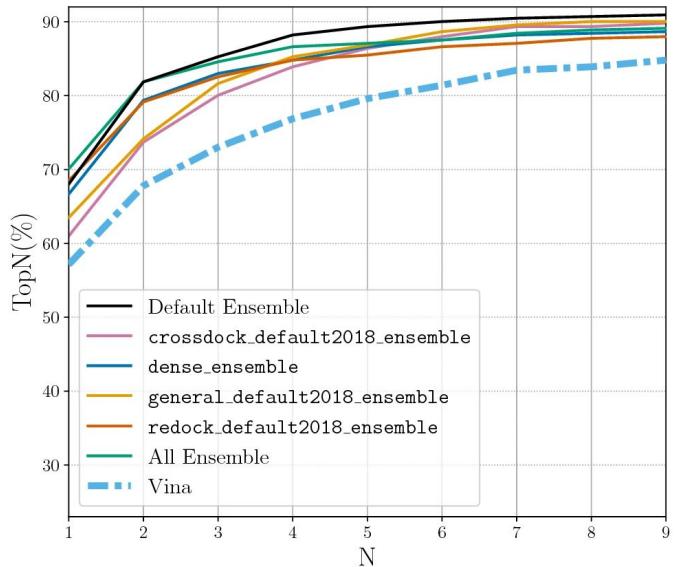


(a) Redocking

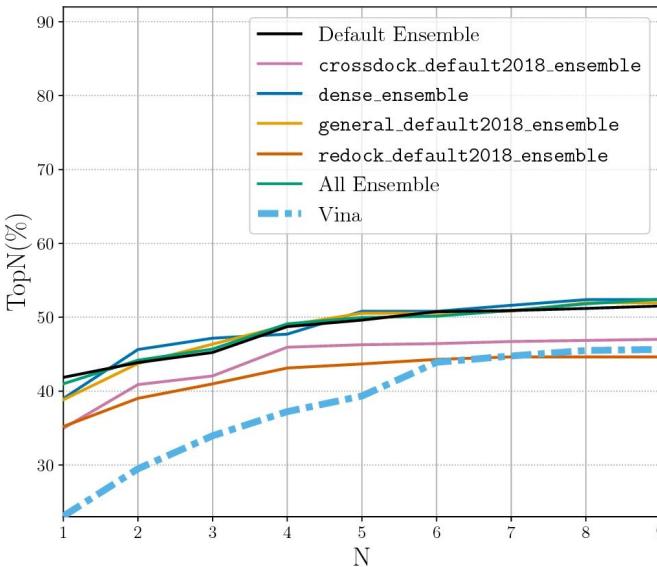


(b) Cross-docking

# AI-powered Molecular Docking



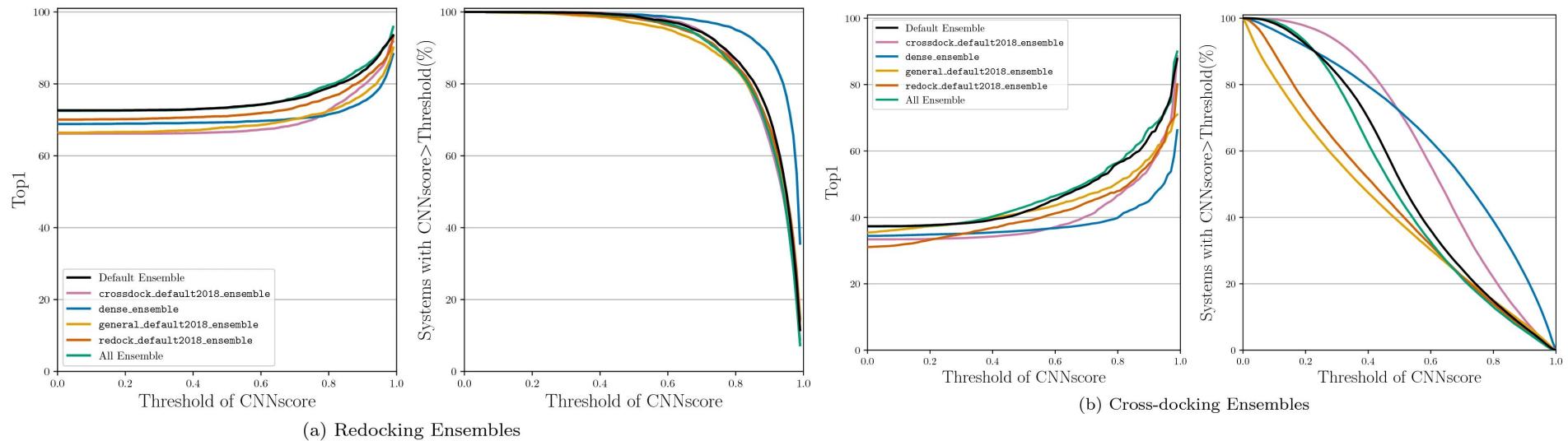
(a) Redocking Ensembles



(b) Cross-docking Ensembles

Testing against data not used in the training set

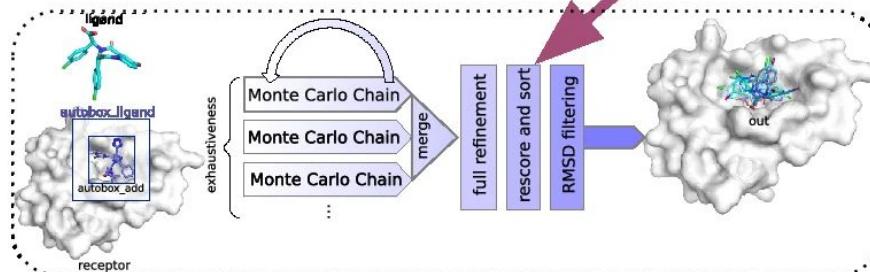
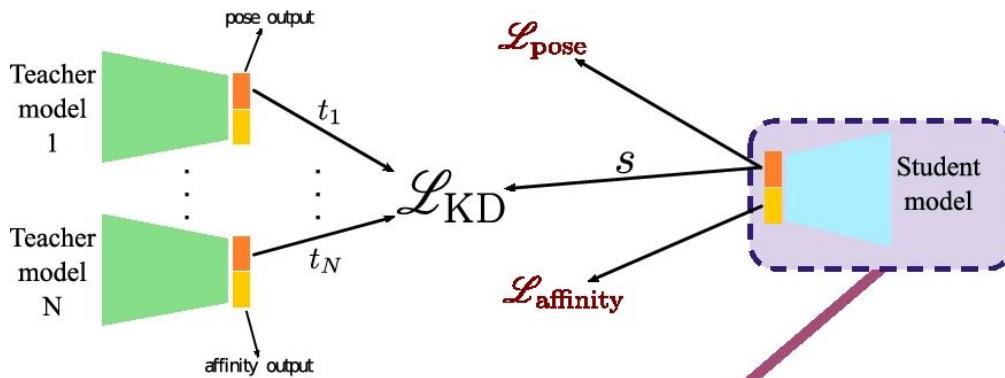
# AI-powered Molecular Docking



# AI-powered Molecular Docking

Gnina 1.3  
updates

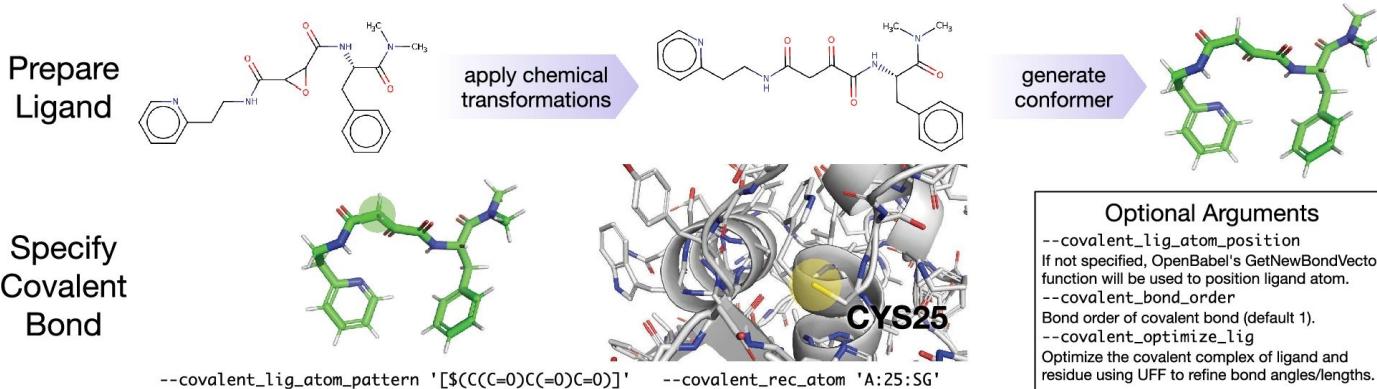
Faster  
calculations



# AI-powered Molecular Docking

Gnina 1.3 updates

Covalent docking  
(atom of the ligand  
is covalently bound  
to an atom of the  
receptor)

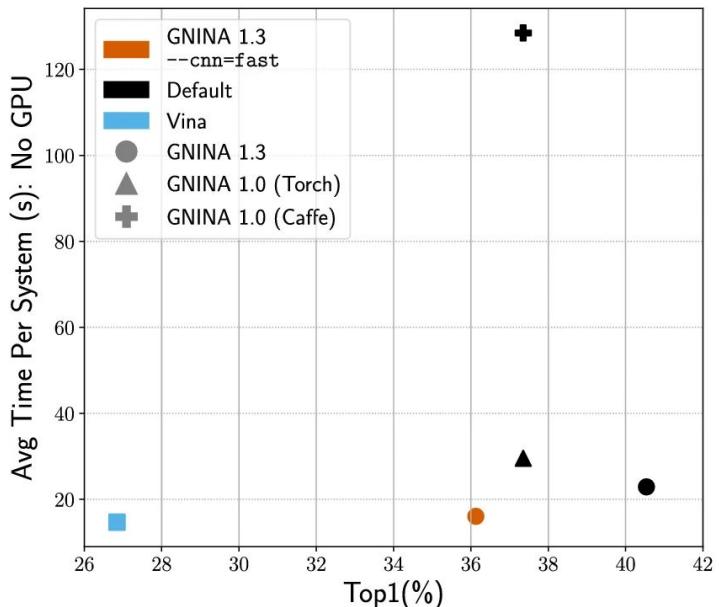


# AI-powered Molecular Docking

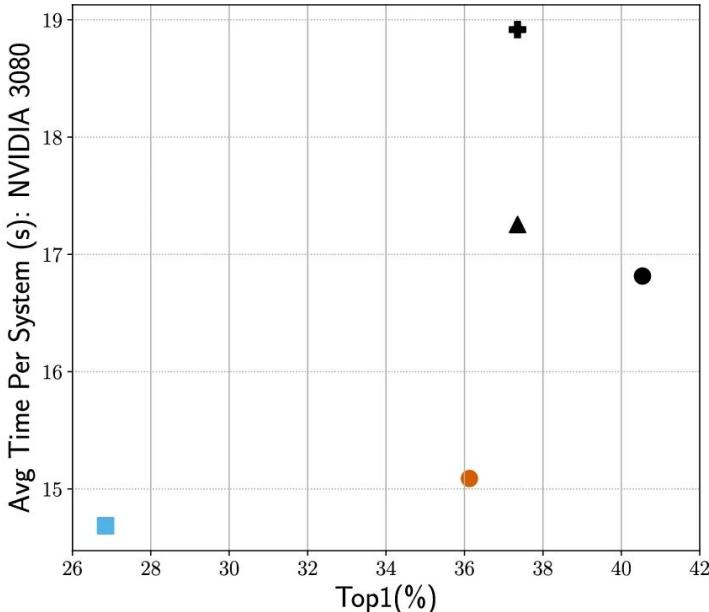
Gnina 1.3 updates

Old default ensemble is almost as fast as Vina now!  
(`--cnn=fast`)

Good for high throughput VS



(a) No GPU

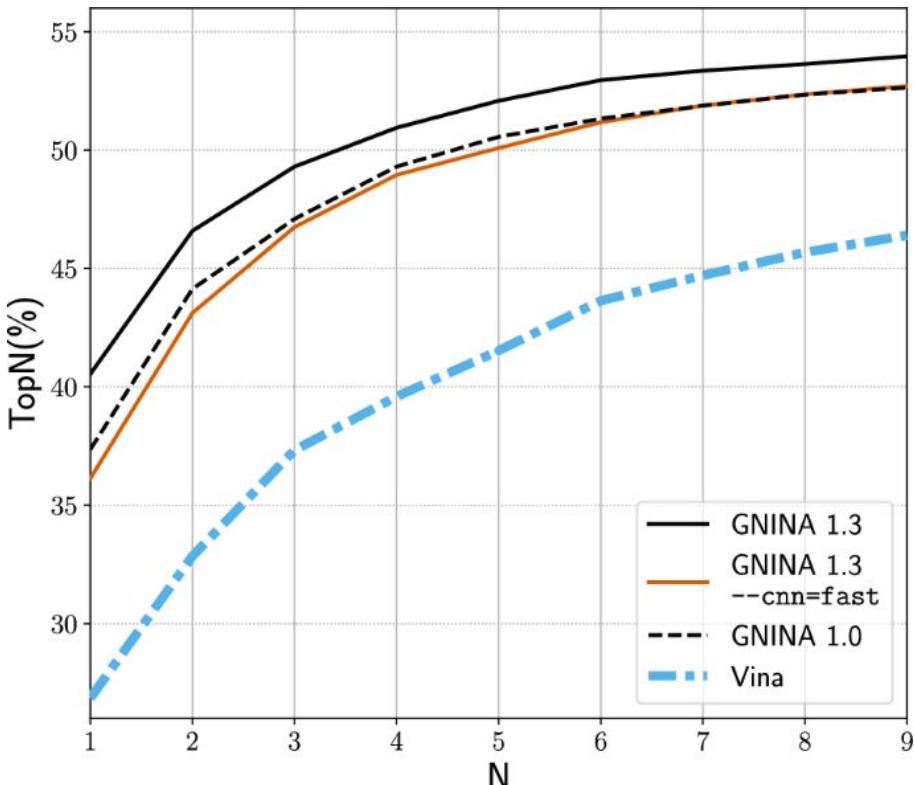


(b) GPU

# AI-powered Molecular Docking

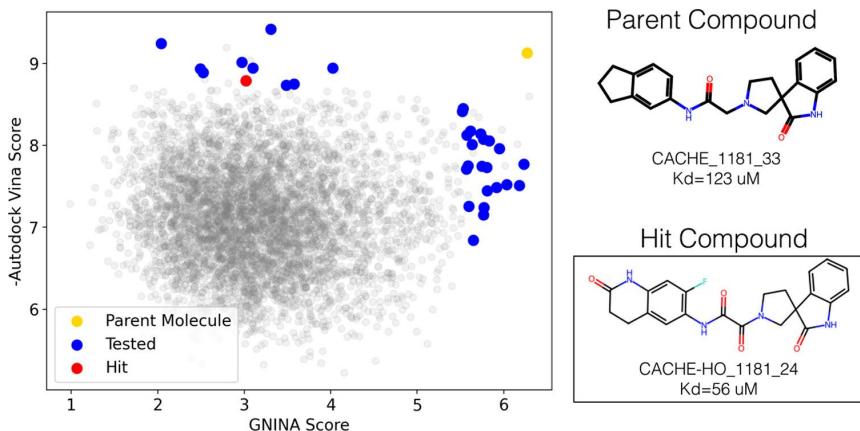
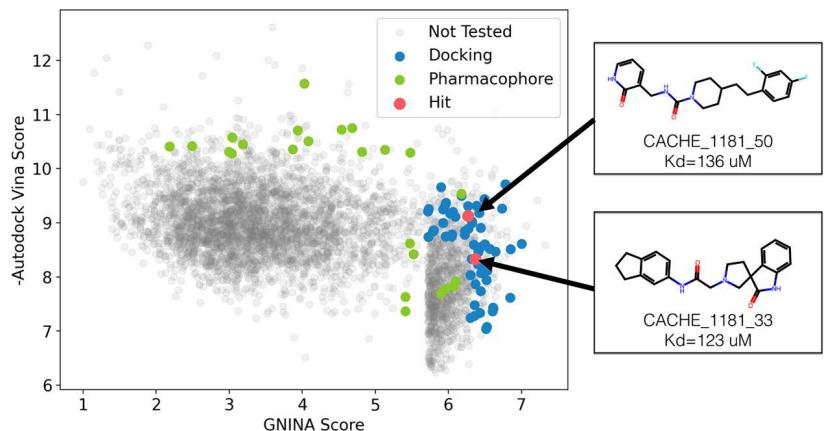
Gnina 1.3 updates

New default ensemble performs better than the rest



# CACHE Challenge #1

Critical Assessment of Computational Hit-finding Experiments: A public–private partnership benchmarking initiative to enable the development of computational methods for hit-finding



Gnina tied in the first place in the Challenge #1!

*“...we find it noteworthy that a purely docking-based approach (with Gnina) was not obviously distinguishable from more complex approaches that used molecular simulation or were guided by expert medicinal chemist.”*

# AI-powered Molecular Docking - Practical Lesson

- We're using Colab notebooks to perform molecular docking using Gnina!
- You can access the notebooks via the link below or the QR code:

<https://github.com/pablo-arantes/EGB2025-MC14>

