

# Supporting Information:

## Scaffold-Constrained Molecular Generation

Maxime Langevin,<sup>†,‡</sup> Hervé Minoux,<sup>‡</sup> Maximilien Levesque,<sup>\*,†,¶</sup> and Marc Bianciotto<sup>\*,‡</sup>

<sup>†</sup>*PASTEUR, Département de chimie, École Normale Supérieure, PSL University, Sorbonne Université, CNRS, 75005 Paris, France*

<sup>‡</sup>*Molecular Design Sciences - Integrated Drug Discovery, Sanofi R&D, 94400 Vitry-sur-Seine, France*

<sup>¶</sup>*Aqemia, 75001 Paris, France*

E-mail: \*maximilien.levesque@aqemia.com; \*marc.bianciotto@sanofi.com

## Data curation and software availability

### Code availability

All software and data to reproduce the results of this paper are available at:

<https://github.com/maxime-langevin/scaffold-constrained-generation>.

To implement the methods presented above, we build on the existing codebase of Olivecrona et al. , available at <https://github.com/MarcusOlivecrona/REINVENT>. The different experiments are reproduced in Jupyter notebooks<sup>S2</sup> available in our codebase. An extra notebook showing basic usage for researchers interested in simply using our method without necessarily reproducing our results is also available.

## ChEMBL dataset

The ChEMBL<sup>S3</sup> database is often used to train generative models of drug-like molecules. To train our RNN, we use a preprocessed version of the ChEMBL<sup>S1</sup> where only molecules having between 10 and 50 heavy atoms and comprised of elements  $\in \{H, B, C, N, O, F, Si, P, S, Cl, Br, I\}$  were kept. The original filtered ChEMBL dataset can be found and downloaded at: [https://github.com/MarcusOlivecrona/REINVENT/blob/master/data/ChEMBL\\_filtered](https://github.com/MarcusOlivecrona/REINVENT/blob/master/data/ChEMBL_filtered), and in our repository at [https://github.com/maxime-langevin/scaffold-constrained-generation/data/ChEMBL\\_filtered](https://github.com/maxime-langevin/scaffold-constrained-generation/data/ChEMBL_filtered). Furthermore, we filtered the dataset to exclude molecules having one the 17 validation scaffolds as a substructure, yielding the final dataset at [https://github.com/maxime-langevin/scaffold-constrained-generation/data/ChEMBL\\_without\\_sureChEMBL.smi](https://github.com/maxime-langevin/scaffold-constrained-generation/data/ChEMBL_without_sureChEMBL.smi).

## SureChEMBL dataset

The SureChEMBL<sup>S4</sup> database is comprised of patented compounds. The database can be downloaded at <https://chembl.gitbook.io/chembl-interface-documentation/downloads>. 34000 compounds were extracted from SureChEMBL v2019.10.01. Compounds were clustered by Bemis-Murcko scaffold<sup>S5</sup> and 18 chemical series (every molecule in each series having the same scaffold) were kept to be used as a validation set. The molecules in the 18 series can be found at [https://github.com/maxime-langevin/scaffold-constrained-generation/data/SureChEMBL/200323\\_SureChemBL\\_dataset\\_636.sdf](https://github.com/maxime-langevin/scaffold-constrained-generation/data/SureChEMBL/200323_SureChemBL_dataset_636.sdf), and the 18 scaffolds at [https://github.com/maxime-langevin/scaffold-constrained-generation/data/SureChEMBL/surechembl\\_scaffolds.sdf](https://github.com/maxime-langevin/scaffold-constrained-generation/data/SureChEMBL/surechembl_scaffolds.sdf).

## DRD2 dataset

The full DRD2 dataset can be found at [https://github.com/undeadpixel/reinvent-scaffold-decorator/blob/master/training\\_sets/drd2.excapedb.smi.gz](https://github.com/undeadpixel/reinvent-scaffold-decorator/blob/master/training_sets/drd2.excapedb.smi.gz). The scaffold

folds used for the goal-directed benchmark are the ones used in Arús-Pous and al Arús-Pous et al. , and can be found in our codebase at [https://github.com/maxime-langevin/scaffold-constrained-generation/data/DRD2/drd2\\_scaffolds.sdf](https://github.com/maxime-langevin/scaffold-constrained-generation/data/DRD2/drd2_scaffolds.sdf).

## **MMP-12 dataset**

The MMP-12 dataset was downloaded from the supplementary materials of Pickett et al.. The dataset can be found at <https://github.com/maxime-langevin/scaffold-constrained-generation/data/MMP12/mmp12.csv>.

## **Implementation details**

### **Policy masking benchmark**

To assess whether a simple policy masking could achieve scaffold constrained generation, we used the following algorithm:

**Result:** SMILES string with scaffold  $s$

**Input:** scaffold  $s = s_1, \dots, s_n$ , number of sampling steps  $k$

initialize  $h_0$  ;

$x_0 = GO$  ;

$t = 1$  ;

**for**  $i \leftarrow 1$  **to**  $n$  **do**

**if**  $s_i$  *not*  $*$  **then**

        Read  $s_i$  and update  $h_{t-1}$  to  $h_t$  ;

$x_t = s_i$  ;

$t = t + 1$  ;

**else**

**for**  $j \leftarrow 1$  **to**  $k$  **do**

            Sample  $x_t$  from  $P(x|h_t)$  and update  $h'_{t-1}$  to  $h'_t$ ;

$t = t + 1$  ;

**end**

**end**

$x_t = EOS$  ;

**Output:**  $x_0, \dots, x_t$

**Algorithm 1:** Simple policy masking

The algorithm masks the policy to read tokens from the scaffold, and samples a fixed number of tokens when encountering the "\*" symbol that represents an open position. Validity of generated SMILES was averaged over completion of 17 diverse scaffolds corresponding to chemical series from SureChEMBL. The choice of  $k$  (the number of sampling steps, set to 8 in our experiments) was arbitrary and made to roughly reflect the expected size of a decoration in a drug-like molecule. We do not expect the choice of the number of sampling steps to impact significantly the proportion of valid generated SMILES.

## Distribution learning benchmarks

To assess distribution learning benchmarks, 10000 molecules were generated for each scaffold.

### Validity

The validity score for a scaffold is the ratio of the number of valid molecules, as defined in the RDKit<sup>S8</sup>, out of all 10000 generated molecules.

### Unicity

Out of the generated valid molecules, the number of unique molecules is computed as the ratio of molecules with distinct canonical SMILES string.

### Physico-chemical properties

All physico-chemical properties were computed using the RDKit. Properties were computed on the valid molecules out of the 10000 generated for each scaffold, and then grouped together. The overall distributions were plotted against the distributions of both the training and the validation set, in order to check that there was no striking dissimilarity between them.

## Predicting DRD2 activity

One of the major point of the DRD2 activity was to benchmark our method against the Reinvent Scaffold Decorator.<sup>S9</sup> Thus, it seems natural to use the same QSAR model. As we weren't able to find this QSAR model within the codebase reproducing the experiments of the article, we used a QSAR model<sup>S1</sup> used in a work from the same group on the DRD2 dataset, and we assumed that the QSAR model used in the two works was the same. In our codebase, the model used for DRD2 activity prediction can be found at <https://github.com/maxime-langevin/scaffold-constrained-generation/data/clf.pkl>

## Predicting MMP-12 activity

To predict activity on the MMP-12 target, the dataset was split into a training and a test set. Then, a random forest regression algorithm (implemented with Scikit-learn)<sup>S10</sup> was fitted on the training set with continuous targets (corresponding to the experimental  $pIC_{50}$ ), and evaluated on the testing set. The evaluation yielded a coefficient of determination  $r^2 = 0.84$ . The QSAR model is accessible at [https://github.com/maxime-langevin/scaffold-constrained-generation/data/MMP12/final\\_activity\\_model.pkl](https://github.com/maxime-langevin/scaffold-constrained-generation/data/MMP12/final_activity_model.pkl), and evaluation on the test set at [https://github.com/maxime-langevin/scaffold-constrained-generation/MMP12\\_experiments.ipynb](https://github.com/maxime-langevin/scaffold-constrained-generation/MMP12_experiments.ipynb).

## Hill climbing procedure

To optimize molecules in goal-oriented benchmarks, a hill-climbing procedure<sup>S11</sup> was used, as it was shown to be overall the best method amongst different generative models. The algorithm can be summarized as a repetition of the following steps:

- Generate 500 molecules
- Score them and keep the top 50 unique molecules
- Perform 10 rounds of log-likelihood maximization with the 50 best molecules

Those steps are repeated 10 times in a row. The code for performing hill-climbing can be found at [https://github.com/maxime-langevin/scaffold-constrained-generation/hill\\_climbing.py](https://github.com/maxime-langevin/scaffold-constrained-generation/hill_climbing.py).

**Data:**  $n_{steps}$ ,  $n_{samples}$ ,  $n_{train}$ ,  $score(m)$ ,  $P_0$

**Result:** Optimized molecules with respect to  $score$

```
for  $t \leftarrow 0$  to  $n_{steps}$  do  
    smiles  $\leftarrow$  SampleSmiles( $P_t$ ,  $n_{samples}$ );  
    training  $\leftarrow$  {};  
    reorder smiles by ascending values of  $score$ ;  
    for  $j \leftarrow 0$  to  $n_{train}$  do  
         $s \leftarrow$  i-th element of smiles;  
        training  $\leftarrow$  training  $\cup$   $s$ ;  
    end  
     $P_{t+1} \leftarrow$  TrainPolicyOnBatch( $P_t$ , training);  
end
```

**Algorithm 2:** Reinforcement learning

## Figures

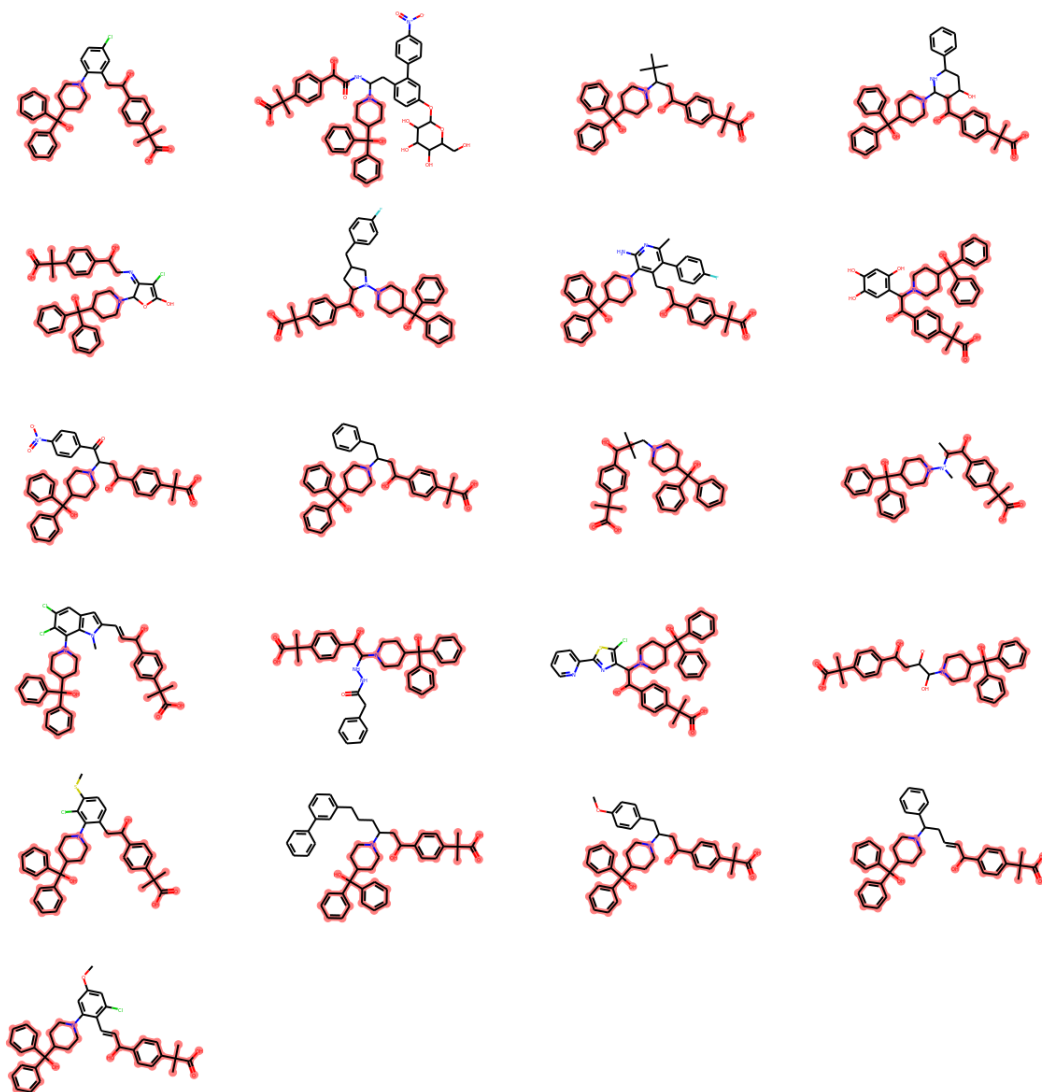


Figure S1: Molecules obtained after modifying the core of fexofenadine. The scaffold constraint is highlighted.

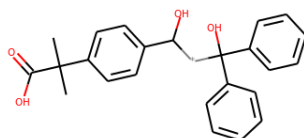


Figure S2: Fexofenadine with one of the rings replaced by an open position



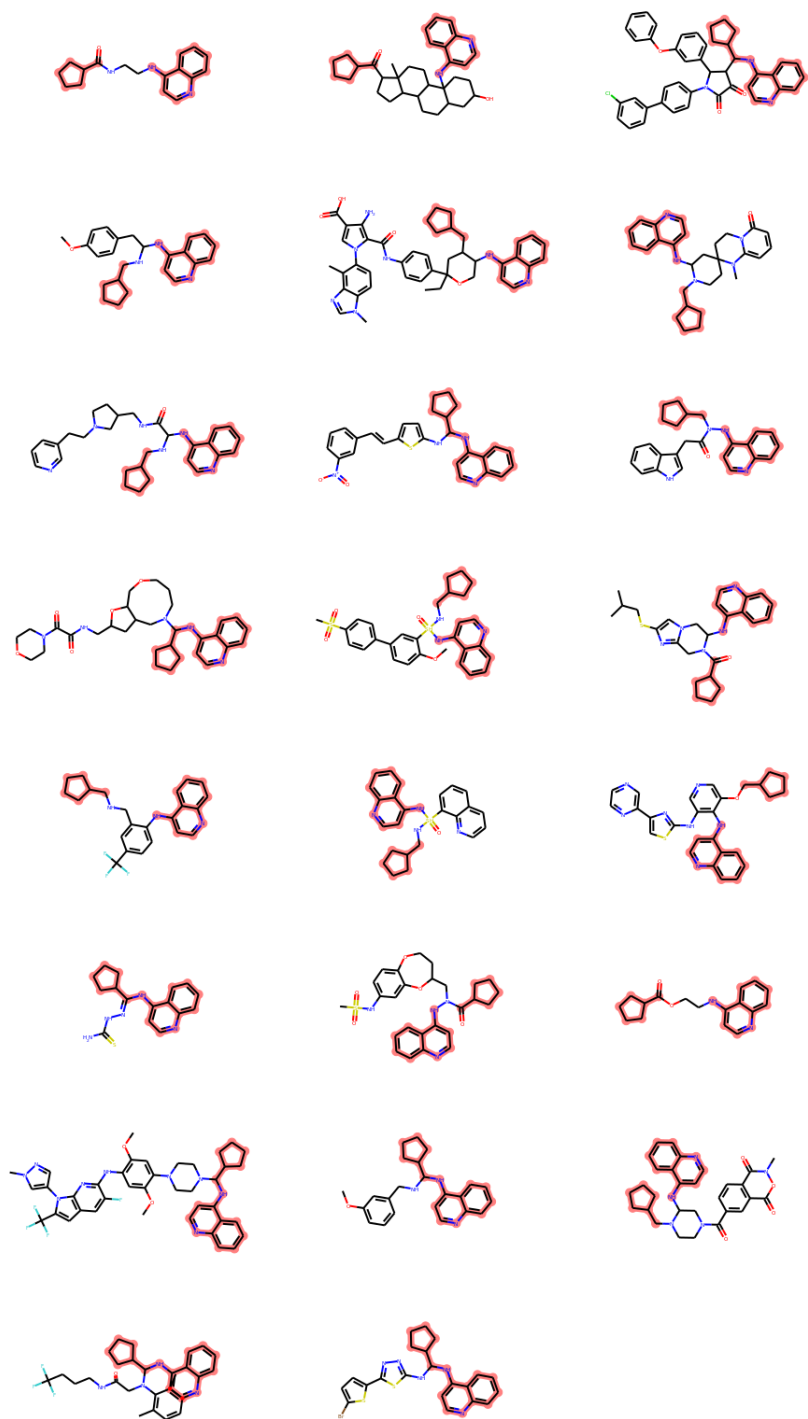


Figure S3: Molecules obtained after modifying the ring system of fexofenadine. The scaffold constraint is highlighted.

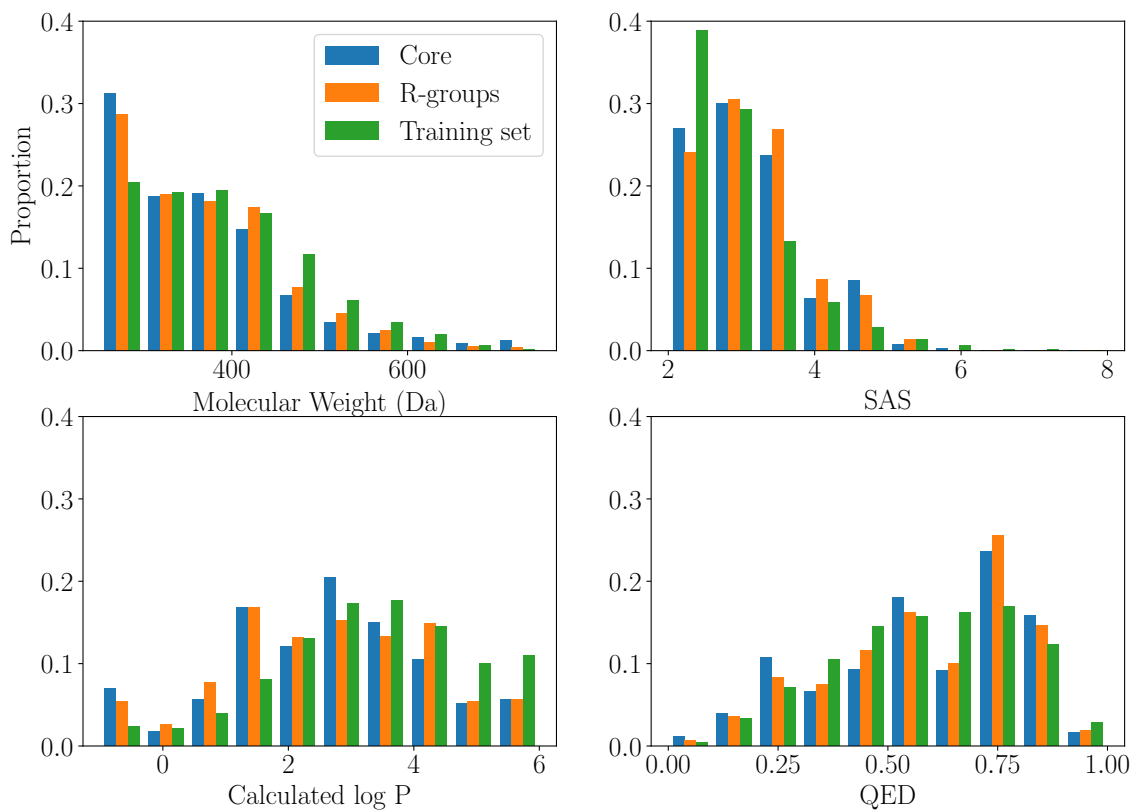


Figure S4: Molecular descriptors distributions for a reference set and molecules sampled around scaffolds with an open position in the core or in a branched decoration

Open positions, either branched or within the core of the molecule, were randomly chosen for each of the 17 scaffolds from the validation set (10 times for each scaffold). Each uncompleted scaffold was then sampled 25 times.

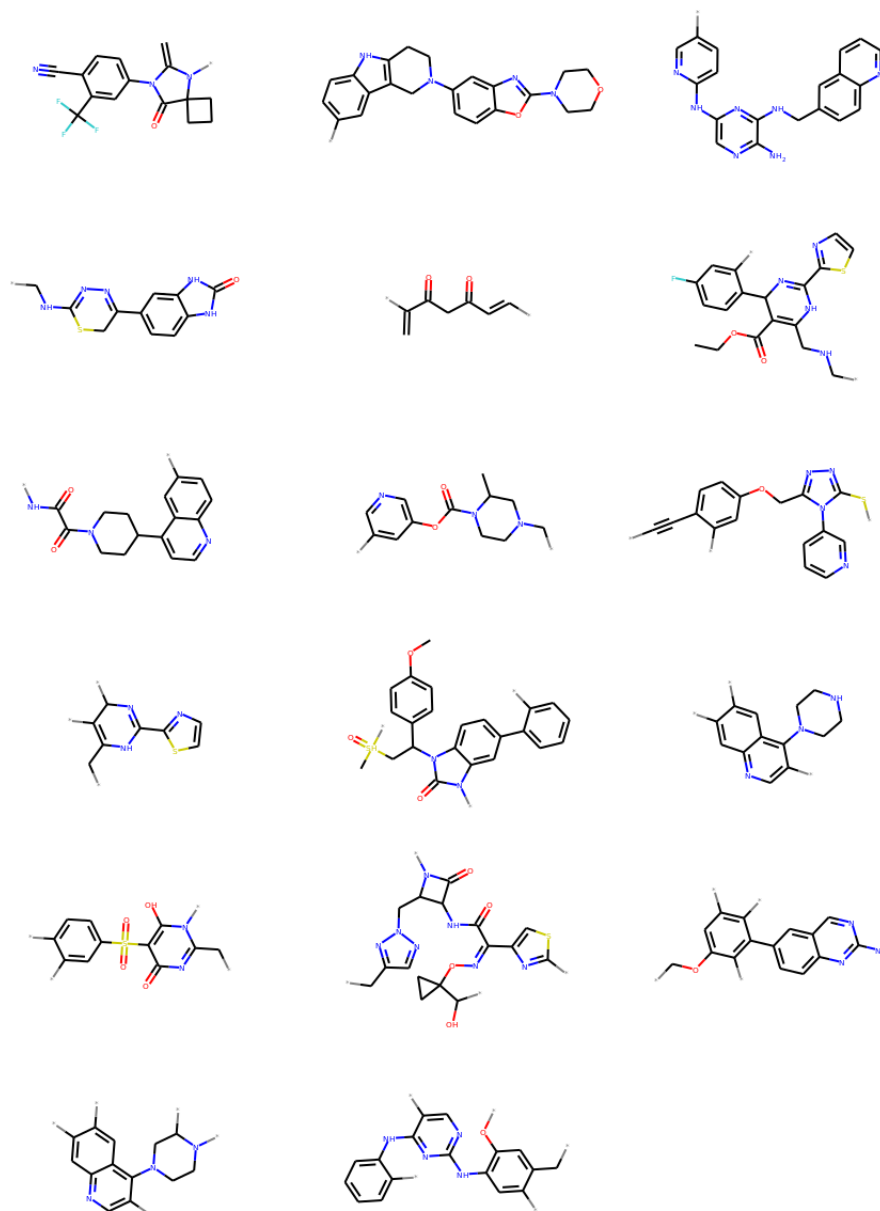


Figure S5: Validation scaffolds from SureChEMBL

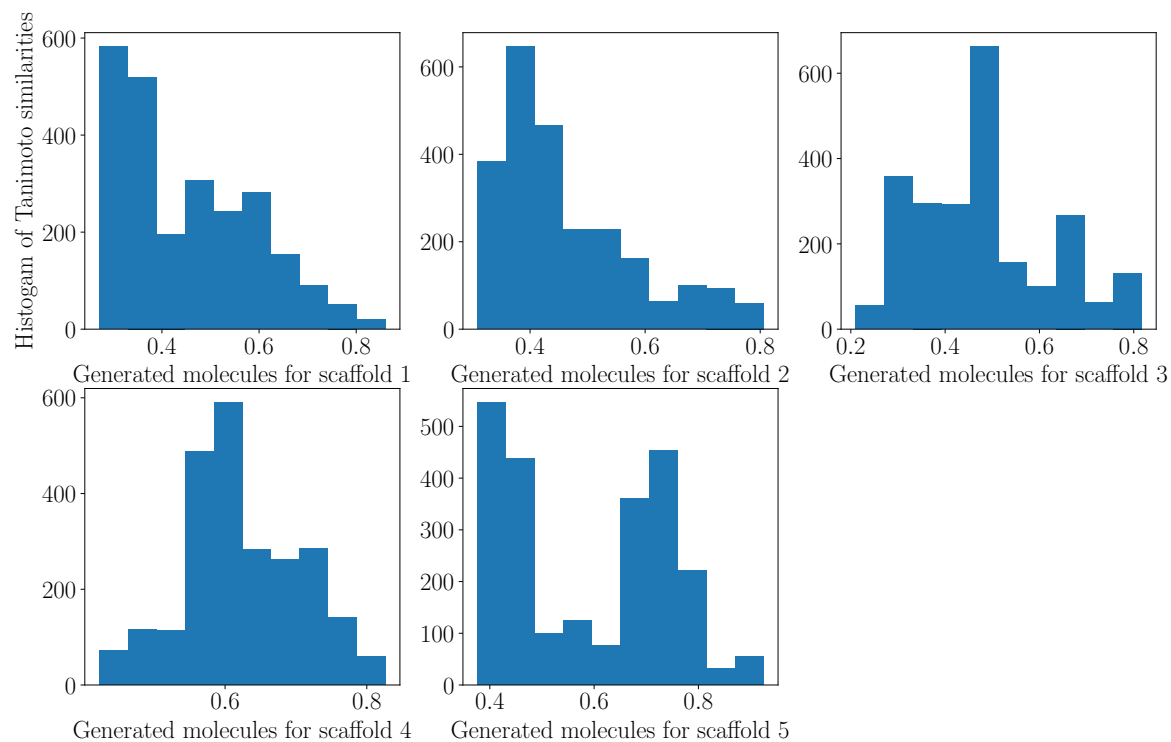


Figure S6: Tanimoto similarity histograms between generated molecules for every one of five DRD2 scaffold

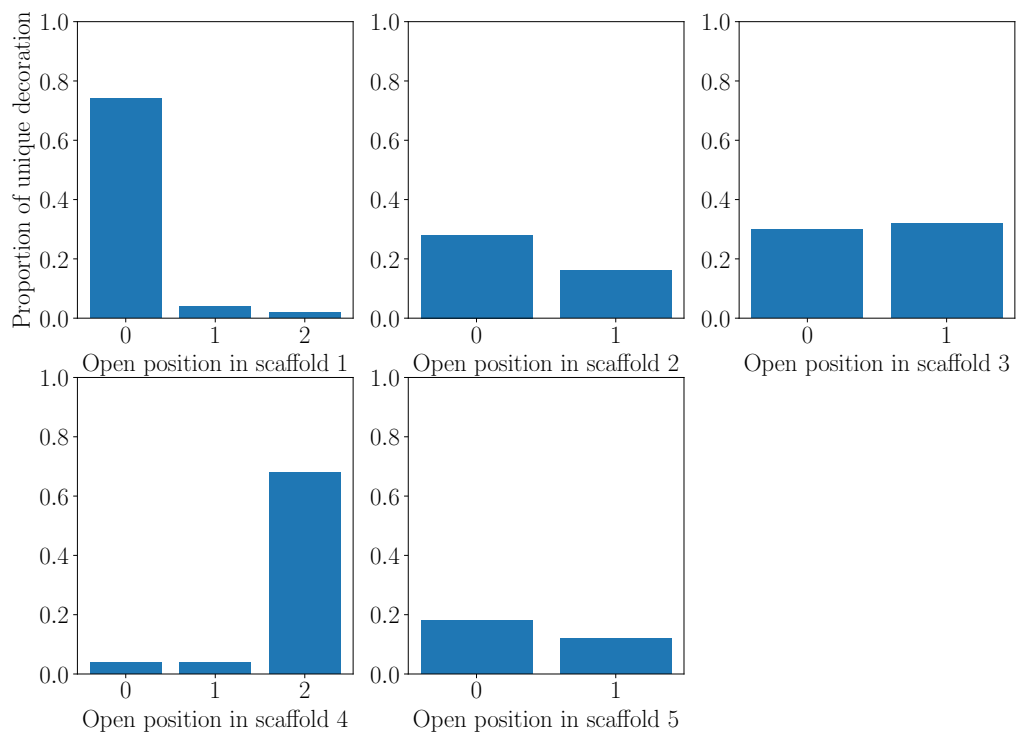


Figure S7: Proportion of unique decoration for each open position in every one of five DRD2 scaffold

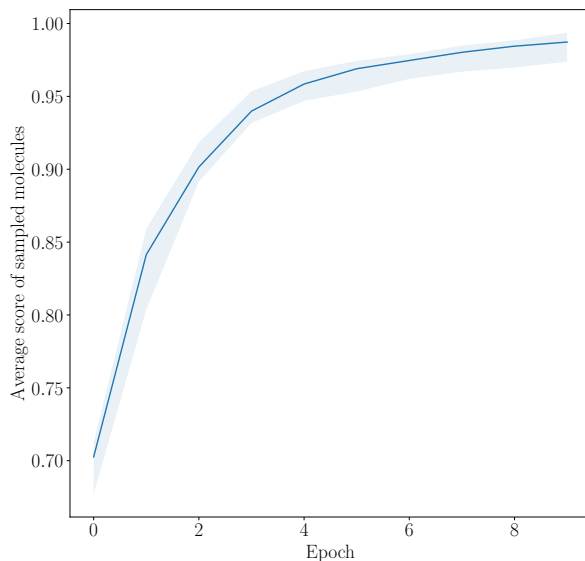


Figure S8: Evolution of average scores of 50 highest scoring molecules generated at each step of hill-climbing, with confidence interval at 95% over 10 runs

## References

- (S1) Olivecrona, M.; Blaschke, T.; Engkvist, O.; Chen, H. Molecular De Novo Design through Deep Reinforcement Learning. *J. Cheminf.* **2017**, *9*, 48.
- (S2) Kluyver, T.; Ragan-Kelley, B.; Pérez, F.; Granger, B.; Bussonnier, M.; Frederic, J.; Kelley, K.; Hamrick, J.; Grout, J.; Corlay, S.; Ivanov, P.; Avila, D.; Abdalla, S.; Willing, C. In *Positioning and Power in Academic Publishing: Players, Agents and Agendas*; Loizides, F., Schmidt, B., Eds.; 2016; pp 87 – 90.
- (S3) Gaulton, A.; Hersey, A.; Nowotka, M.; Bento, A. P.; Chambers, J.; Mendez, D.; Atkinson, F.; Mutowo, P.; Bellis, L. J.; Cibrián-Uhalte, E.; Davies, M.; Dedman, N.; Karlsson, A.; Magariños, M. P.; Overington, J. P.; Papadatos, G.; Smit, I.; Leach, R., A. The ChEMBL database in 2017. *Nucleic Acids Res.* **2016**, *45*, D945–D954.
- (S4) Papadatos, G.; Davies, M.; Dedman, N.; Chambers, J.; Gaulton, A.; Siddle, J.;

- Koks, R.; Irvine, S. A.; Pettersson, J.; Goncharoff, N.; Hersey, A.; Overington, J. P. SureChEMBL: a large-scale, chemically annotated patent document database. *Nucleic Acids Res.* **2015**, *44*, D1220–D1228.
- (S5) Bemis, G. W.; Murcko, M. A. The Properties of Known Drugs. 1. Molecular Frameworks. *J. Med. Chem.* **1996**, *39*, 2887–2893.
- (S6) Arús-Pous, J.; Patronov, A.; Bjerrum, E. J.; Tyrchan, C.; Reymond, J.-L.; Chen, H.; Engkvist, O. SMILES-based deep generative scaffold decorator for de-novo drug design. *J. Cheminf.* **2020**, *12*, 38.
- (S7) Pickett, S. D.; Green, D. V. S.; Hunt, D. L.; Pardoe, D. A.; Hughes, I. Automated Lead Optimization of MMP-12 Inhibitors Using a Genetic Algorithm. *ACS Med. Chem. Lett.* **2010**, *2*, 28–33.
- (S8) Landrum, G. RDKit: Open-source cheminformatics. 2020; <http://www.rdkit.org>.
- (S9) Arús-Pous, J.; Johansson, S. V.; Prykhodko, O.; Bjerrum, E. J.; Tyrchan, C.; Reymond, J.-L.; Chen, H.; Engkvist, O. Randomized SMILES strings improve the quality of molecular generative models. *J. Cheminf.* **2019**, *11*, 71.
- (S10) Pedregosa, F.; Varoquaux, G.; Gramfort, A.; Michel, V.; Thirion, B.; Grisel, O.; Blondel, M.; Prettenhofer, P.; Weiss, R.; Dubourg, V.; Vanderplas, J.; Passos, A.; Cournapeau, D.; Brucher, M.; Perrot, M.; Duchesnay, E. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research* **2011**, *12*, 2825–2830.
- (S11) Brown, N.; Fiscato, M.; Segler, M. H.; Vaucher, A. C. GuacaMol: Benchmarking Models for de Novo Molecular Design. *J. Chem. Inf. Model.* **2019**, *59*, 1096–1108.