Supplementary Information (SI) for Chemical Science. This journal is © The Royal Society of Chemistry 2025

**Supplementary Information** *for the article* 

# Digitization of molecular complexity with machine learning

Andrei S. Tyrin,<sup>#</sup> Daniil A. Boiko,<sup>#</sup> Nikita I. Kolomoets, Valentine P. Ananikov\*

Zelinsky Institute of Organic Chemistry, Russian Academy of Sciences, Leninsky prospekt 47, Moscow, 119991, Russia; http://AnanikovLab.ru \*e-mail: val@ioc.ac.ru; <sup>#</sup>equal contributions

#### Contents

| Dataset  | S2  |
|--|---|
| Active Learning                                  | S3  |
| Model Selection                                  | S4  |
| Interpreting molecular complexity                | S6  |
| FDA drugs results                                | S7  |
| Total syntheses with molecular complexity values | S10   |
| Datasets used in the analysis                    | S16   |
| ferences   | S16   |
| f  | Dataset<br>Active Learning<br>Model Selection<br>Interpreting molecular complexity<br>FDA drugs results<br>Total syntheses with molecular complexity values<br>Datasets used in the analysis<br>erences |

### 1. Dataset

As was mentioned, the collected dataset consisted of three phases (consult the paper for detailed description and motivation for each of them). Using the quality controls and group membership as filtering criteria, we obtained the following count distribution for each of the phases:

- 1) Absolute values phase 10978 labels used.
- 2) Active learning phase 94810 labels used.
- 3) Ranking phase 90295 labels used.



The sizes of three phases compared in the Figure S1.

Figure S1. Distribution of votes number across different phases.

Active learning allowed to collect labels for 164017 molecules, which is a significant increase compared to other data-based approaches tackling the molecular complexity determination.

Each expert was asked their education level. The results are presented in Figure S2. As can be seen, the largest group of experts consisted of PhD students pursuing organic chemistry as their domain of expertise followed by experts who already obtained PhD in organic chemistry.



**Figure S2**. Education level among the experts whose labels served as labels for the ranking model training.

# 2. Active Learning

After collecting the initial set of labels with randomly sampled groups of molecules the active learning was employed [links]. The reason for integrating active learning into the data collection pipeline is to cover as much of the chemical space by using as little of human effort as possible. During the most productive days of data collection, the ensemble of ranking models was retrained nightly using the already collected data, after which it was used to select a batch of unlabeled molecules using uncertainty in the ranking score as a criterion. After this step, selected molecules were grouped into 5 clusters (corresponding to 5 complexity thresholds), and within each cluster most similar

molecules were grouped into samples that were provided to the experts on the following day. This step, thus achieved two goals: covering the chemical space as efficient as possible as well as collecting the labels for the molecules that the model considers highly similar.

#### 3. Model Selection

Although having different data collection phases is beneficial simple merging of them for the subsequent training might be hard due to the different sizes and types of molecules in each phase. For instance, the second phase is based on the data collected during the first phase and data collected during active learning. Therefore, we used the group weight parameter between different phases of the GBDTs model as a hyperparameter for the selection of best model. Other parameters included:

"loss\_function": "YetiRank", "iterations": 1000, "depth": 8, "learning\_rate": 0.05, "bootstrap\_type": "MVS",

We found that changing them didn't significantly change the performance of the model on the test dataset that was obtained by random split. Therefore, the main hyperparameter that was used corresponded to the group weight. The weight was discretized up to 0.1 thus resulting in 66 different models. Pair Accuracy and Function Group Test (FGT) was used as selection criteria. The functional groups included: phenyl, methyl, amine, hydroxyl, carboxyl, aldehyde, t-Bu, isopropyl, nitril, sulfo, fluorine, chlorine, bromine, iodine, and methoxy groups.

By equipping this approach for the hyperparameter tuning we were able to combine the data collected in different phases by maximizing the ranking quality as well as having the meaningful MC predictions validated by FGT. The detailed results of the conducted analysis are presented in the Figure S3. As was mentioned in the main part of the article, the selected phase weights correspond to 70%, 10%, and 20% distributed

between the first, second, and third phases respectively. This distribution corresponds to the optimal joint performance from both PA and FGT perspectives.



**Figure S3**. Ternary diagram visualization for the performance metrics. (a) pair accuracy on different test phases obtained by random train/test split; (b) averaged pair accuracy for different phases; (c) FGT results (the higher score/darker color shade is better).

# 4. Interpreting molecular complexity

In addition to Fig. 1a, the complete SHAP beeswarm plot is given in Fig. S3.





To further interpret the relationship between the molecular complexity and other molecular features that were used for the training of the machine learning model, Fig. S5 was made.



**Figure S5**. MC plotted against the (a) SCScore, (b) molecular weight, (c) number of aromatic rings, and (d) topological polar surface area. Molecules were taken from the ChEMBL database.

Fig. S5 allows to explore further what guides scientists during the analysis of molecules from the MC perspective.

# 5. FDA drugs results

Additionally to the 2D histogram presented in the Fig. 3a of the main text, we provide a 3D histogram (Fig. S6) to provide an additional perspective on the molecular complexity evolution of small molecule drugs. Additionally, Fig. S6 contains the results of the Mann-Kendall test, indicating clear growth of the molecular complexity over years.



**Figure S6**. (a) 3D histogram illustrating the MC distribution of drugs over time; (b) histogram with the results of Mann-Kendall test;

In addition to the molecules of the median molecular complexity per 5 years intervals illustrated in the main text, we provide the MC-annotated molecular structures in Fig. S7.



1985: nabilone, MC=4.3 1986: permethrin, MC=4.76 1987: lovastatin, MC=5.76 1988: nizatidine, MC=3.94 1989: omeprazole, MC=4.85

1990: isradipine, MC=5.0 1991: loracarbef, MC=5.0 1992: atovaquone, MC=4.74 1993: fenofibrate, MC=4.16 1994: famiciclovir, MC=4.62



1995: nisoldipine, MC=4.66 1996: topiramate, MC=5.23 1997: tazarotene, MC=4.35 1998: celecoxib, MC=5.19 1999: rofecoxib, MC=4.27



2000: meloxicam, MC=4.9 2014: eliglustat, MC=5.04 2015: cholic acid, MC=6.13 2021: infigratinib, MC=6.06 2001: travoprost, MC=5.81



2002: ezetimibe, MC=5.5 2003: aprepitant, MC=6.34 2004: eszopidone, MC=5.46 2005: deferasirox, MC=5.11 2006: ranolazine, MC=5.25

your where the

of the

2007: aliskiren, MC=5.72 2008: clevidipine, MC=5.0 2009: besifloxacin, MC=5.37 2010: ulipristal acetate, MC=5.97 2011: abiraterone acetate, MC=5.47

and

5

for many man

2012: bosutinib, MC=5.98 2013: alogliptin, MC=4.84 2022: daridorexant, MC=5.97 2016: pimavanserin, MC=5.11 2017: enasidenib, MC=5.95

afre soll

2018: binimetinib, MC=5.62 2019: darolutamide, MC=5.64 2020: remimazolam, MC=5.55

Figure S7. FDA approved drugs with closest to the average molecular complexity value per each year.

# 6. Total syntheses with molecular complexity values

The schemes below contain the synthetic route strychnine and artemisinin. All the schemes contain the intermediates with the biggest molecular complexity values.

#### 1. Strychnine

a. Woodward 1954<sup>1</sup>



# b. Overman 1993<sup>2</sup>



6.37

S11





5.95

d. Vanderwal 2011<sup>4</sup>



e. Biosynthesis<sup>5</sup>



# 2. Artemisinin





# 7. Datasets used in the analysis

Several datasets were used during the performed analysis. In particular,

- Figure 1b relies on the sample of PubChem that was obtained according to the description given in the Methods section.
- Results given in Figure 2a-b are the selection of ChEMBL drugs that can be reproduced by using the Datamol framework (the dataset itself can be accessed by datamol.data.chembl\_drugs(as\_df=True)). Figure 2c involves the analysis of common benchmarking datasets. In particular, QM9, Tox21, and HIV datasets were considered. The latter two were accessed through the therapeutics data commons API.<sup>10</sup>
- Reaction Atlas in Figure 3 is constructed based on the Schneider 50k<sup>11</sup> dataset and is processed according to the methodology developed by Schwaller<sup>12</sup> et al.
- Figure 4a relies on the collection of FDA-approved small molecule drugs, which are publicly available. Additionally, the datasets collected by Ross et al. was used for the quantitative analysis presented in Figure 4b-d. This involves both measurements (that were received as described in corresponding publication<sup>13</sup>)as well as docked structures.
- Results in Figure 5 involve the synthetic routes of two natural products that were manually drawn and parsed by the authors of this work.

# References

- Woodward, R. B.; Cava, M. P.; Ollis, W. D.; Hunger, A.; Daeniker, H. U.; Schenker, K. THE TOTAL SYNTHESIS OF STRYCHNINE. *J Am Chem Soc* 1954, 76 (18), 4749–4751. https://doi.org/10.1021/ja01647a088.
- (2) Knight, S. D.; Overman, L. E.; Pairaudeau, G. Asymmetric Total Syntheses of (-)and (+)-Strychnine and the Wieland-Gumlich Aldehyde. *J Am Chem Soc* **1995**, *117* (21), 5776–5788. https://doi.org/10.1021/ja00126a017.
- (3) Nakanishi, M.; Mori, M. Total Synthesis of (-)-Strychnine. *Angewandte Chemie International Edition* **2002**, *41* (11), 1934. https://doi.org/10.1002/1521-3773(20020603)41:11<1934::AID-ANIE1934>3.0.CO;2-F.
- (4) Martin, D. B. C.; Vanderwal, C. D. A Synthesis of Strychnine by a Longest Linear Sequence of Six Steps. *Chem Sci* 2011, *2* (4), 649. https://doi.org/10.1039/c1sc00009h.
- (5) Hong, B.; Grzech, D.; Caputi, L.; Sonawane, P.; López, C. E. R.; Kamileen, M. O.; Hernández Lozada, N. J.; Grabe, V.; O'Connor, S. E. Biosynthesis of Strychnine. *Nature* **2022**, *607* (7919), *617–622*. https://doi.org/10.1038/s41586-022-04950-4.
- (6) Avery, M. A.; Chong, W. K. M.; Jennings-White, C. Stereoselective Total Synthesis of (+)-Artemisinin, the Antimalarial Constituent of Artemisia Annua L. *J Am Chem Soc* **1992**, *114* (3), 974–979. https://doi.org/10.1021/ja00029a028.
- (7) Zhu, C.; Cook, S. P. A Concise Synthesis of (+)-Artemisinin. *J Am Chem Soc* **2012**, *134* (33), 13577–13579. https://doi.org/10.1021/ja3061479.
- (8) Krieger, J.; Smeilus, T.; Kaiser, M.; Seo, E.; Efferth, T.; Giannis, A. Total Synthesis and Biological Investigation of (−) Artemisinin: The Antimalarial Activity of

Artemisinin Is Not Stereospecific. *Angewandte Chemie International Edition* **2018**, *57* (27), 8293–8296. https://doi.org/10.1002/anie.201802015.

- (9) Ikram, N. K. B. K.; Simonsen, H. T. A Review of Biotechnological Artemisinin Production in Plants. *Front Plant Sci* **2017**, *8*. https://doi.org/10.3389/fpls.2017.01966.
- Huang, K.; Fu, T.; Gao, W.; Zhao, Y.; Roohani, Y.; Leskovec, J.; Coley, C. W.; Xiao, C.; Sun, J.; Zitnik, M. Artificial Intelligence Foundation for Therapeutic Science. *Nat Chem Biol* 2022, *18* (10), 1033–1036. https://doi.org/10.1038/s41589-022-01131-2.
- (11) Schneider, N.; Lowe, D. M.; Sayle, R. A.; Landrum, G. A. Development of a Novel Fingerprint for Chemical Reactions and Its Application to Large-Scale Reaction Classification and Similarity. *J Chem Inf Model* **2015**, *55* (1), 39–53. https://doi.org/10.1021/ci5006614.
- (12) Schwaller, P.; Probst, D.; Vaucher, A. C.; Nair, V. H.; Kreutter, D.; Laino, T.; Reymond, J. L. Mapping the Space of Chemical Reactions Using Attention-Based Neural Networks. *Nat Mach Intell* **2021**, *3* (2), 144–152. https://doi.org/10.1038/s42256-020-00284-w.
- (13) Shi, L.; Pan, H.; Liu, Z.; Xie, J.; Han, W. Roles of PFKFB3 in Cancer. Signal Transduct Target Ther 2017, 2 (1), 17044. https://doi.org/10.1038/sigtrans.2017.44.