RESEARCH ARTICLE



Deep generative models for ligand-based de novo design applied to multi-parametric optimization

Quentin Perron¹ | Olivier Mirguet^{2,3} | Hamza Tajmouati¹ | Adam Skiredj¹ Anne Rojas^{2,3} | Arnaud Gohier^{2,3} | Pierre Ducrot^{2,3} | Marie-Pierre Bourguignon^{2,3} | Patricia Sansilvestri-Morel^{2,3} | Nicolas Do Huu¹ | Françoise Gellibert^{2,3} Yann Gaston-Mathé¹

Correspondence

Yann Gaston-Mathé, Iktos, 65 rue de Prony, 75017 Paris, France.

Email: yann.gaston.mathe@iktos.com

Abstract

Multi-parameter optimization (MPO) is a major challenge in new chemical entity (NCE) drug discovery. Recently, promising results were reported for deep learning generative models applied to de novo molecular design, but, to our knowledge, until now no report was made of the value of this new technology for addressing MPO in an actual drug discovery project. In this study, we demonstrate the benefit of applying AI technology in a real drug discovery project. We evaluate the potential of a ligand-based de novo design technology using deep learning generative models to accelerate the obtention of lead compounds meeting 11 different biological activity objectives simultaneously. Using the initial dataset of the project, we built QSAR models for all the 11 objectives, with moderate to high performance (precision between 0.67 and 1.0 on an independent test set). Our DL-based AI de novo design algorithm, combined with the QSAR models, generated 150 virtual compounds predicted as active on all objectives. Eleven were synthetized and tested. The Aldesigned compounds met 9.5 objectives on average (i.e., 86% success rate) versus 6.4 (i.e., 58% success rate) for the initial molecules measured on all objectives. One of the AI-designed molecules was active on all 11 measured objectives, and two were active on 10 objectives while being in the error margin of the assay for the last one. The AI algorithm designed compounds with functional groups, which, although being rare or absent in the initial dataset, turned out to be highly beneficial for the MPO.

KEYWORDS

artificial intelligence, drug discovery, lead-optimization, multiparameter optimization

INTRODUCTION 1

Drug design is a challenging task. From hit identification to hit-to-lead and lead optimization, the quest to discover a new chemical entity (NCE) with desired properties is burdensome. Exploration of a nearly infinite chemical space (10⁶⁰ drug-like molecules is a low range figure)^{1,2} is required in order to solve a multi-parametric optimization (MPO) challenge: identifying the rare compounds which satisfy all the objectives of the project, such as biological activity, selectivity, (lack of) toxicity, pharmacokinetics (i.e., DMPK), synthetic accessibility and finally novelty.^{3,4} The average cost to develop a pre-clinically validated drug candidate is estimated around \$50 million, and drug design, more specifically lead optimization, represents the lion's share $(\sim 70\%)$ of the cost of preclinical research.⁵

Structure- and ligand-based computer aided drug design (CADD) technologies (e.g. docking, QSAR, etc.), which have been

This is an open access article under the terms of the Creative Commons Attribution-NonCommercial-NoDerivs License, which permits use and distribution in any medium, provided the original work is properly cited, the use is non-commercial and no modifications or adaptations are made. © 2022 The Authors. Journal of Computational Chemistry published by Wiley Periodicals LLC.

J Comput Chem. 2022;43:692-703. 692 wileyonlinelibrary.com/journal/jcc

¹Iktos, Paris, France

²Institut De Recherches Servier, Suresnes, France

³Institut De Recherches Servier, Croissy. France

096987x, 2022, 10, Downloaded from https://onlinelibrary.wiley.com/doi/10.1002/jcc.26826 by Cochrane France, Wiley Online Library on [15/05/2025]. See the Terms and Conditions on Wiley Online Library for rules of use; OA articles are governed by the applicable Creative Commons License

developed to improve the productivity of the drug design process, have brought notable progress over the last decades. ^{6,7} Still, most classical CADD approaches have focused on the prediction of molecular properties rather than on the exploration of the chemical space to identify novel compounds with optimal properties. Such in silico exploration of the chemical space has mostly been performed through the virtual screening of pre-existing or virtual compound libraries, with the exploration being intrinsically restricted to the initial compound library. ^{8,9} Graph-based genetic algorithms, sometimes used for in silico chemical optimization, have had limited success and are mostly limited by transformation rules. ¹⁰ More recently, the development of artificial intelligence (AI) approaches to drug discovery, and more specifically de novo drug design through the use of deep generative models, has triggered a lot of interest in the CADD community. ¹¹

Generative models for molecular design can be characterized by three main features: (1) which molecular representation they use; (2) how they generate molecules; and (3) how they perform property optimization. Many methods have been reported, each with different approaches regarding those features: (i) The molecular representation can be either text (SMILES, 12,13 SELFIES, 14,15 DeepSMILES¹⁶), a graph or a set of fragments. 17-19 (ii) The generation strategy can use a simple policy, for instance: add or remove atoms or bonds. 17 It can also rely on deep generative models such as recurrent neural networks (RNNs), auto-encoders (AEs) or generative adversarial networks (GANs). 20,21 (iii) The property optimization strategy can be based on reinforcement learning, 17,22,23 continuous optimization, 20 Bayesian optimization, 24 genetic algorithms 15 or particle swarm optimization.

Despite the amount of research in generative modeling and its potential to allow an efficient exploration of the chemical spaces to identify new molecules with the desired in silico properties, evidence of the benefit of such Al-based approaches to solve MPO issues in complex real-life cases is still elusive, and Al-based drug design is perceived as overhyped by a significant part of the chemists' and chemo-informaticians' community.²⁶

As previously stated, MPO is a major challenge in NCE drug discovery projects, and the inability to identify molecules meeting the Target Product Profile (TPP) in LO is an important cause of NCE project failure or delay. Some recent works were conducted in order to generate new molecules in MPO projects, leading to interesting results, 25,27 however, none of them used real project datasets. Herein, we describe the application of a ligand-based de novo design AI technology based on deep generative models in a real-life LO stage drug discovery project and its impact on fostering the discovery of optimized lead compounds meeting the project's TPP criteria. This study was conducted in 2017 and used a Long Short-Term Memory (LSTM) neural network trained on ChEMBL using teacher forcing with a multi-objective reward function. Since then, works from many research groups have led to the development of more sophisticated generative AI methods for drug design, however this work provides evidence of prospective real-life validation of this technology.

2 | METHODS, DATA AND SOFTWARE STATEMENT

2.1 | Project dataset

The dataset was provided by Servier from an internal and real drug discovery project at LO stage that had been running for several years. The project dataset consisted in a library of 881 molecules with associated bioactivity measurements from 11 biological assays: one primary activity assay (undisclosed phenotypic assay: % of activation at 30 nM), 6 off-target activity assays (selectivity criteria on 5-HT2A, 5-HT2B, alpha1, D1, Na_v1.2, hERG: % of inhibition) and 4 ADME assays (microsomal stability on human (HLM) and rat (RLM): % of stability; permeability and efflux Caco2 assays: % of absorption and efflux ratio). For each objective, a threshold value was defined according to the Target Product Profile (TPP) designed by the project team. A summary of the thresholds, percentage of compounds measured and percentage of compounds meeting the required threshold for each assay is reported in Table 1.

The best molecule from the initial dataset and the 11 Algenerated molecules synthesized and tested are provided as SMILES in the Supplementary material.

2.2 | Software availability

The following software packages were used to perform this work: (1) The QSAR models were built using Scikit-learn;²⁸ (2) Hyperopt was used to optimize the hyperparameters for model selection;²⁹ (3) Training and optimization of the LTSM was performed using Tensorflow;³⁰ (4) Rdkit was used to prepare SMILES, calculate similarities, fingerprints and descriptors.³¹ All the software packages are freely available.

2.3 | QSAR models development

Bioactivity data were binned according to TPP thresholds (i.e., 1 if meeting the TPP specification, else 0). Eleven independent QSAR models were developed using ridge logistic regression based on Morgan fingerprint molecular representations (2048 bits and radius 3).³² The Morgan fingerprint was built without including chirality (two stereoisomers have an identical fingerprint) as most of the molecules in the dataset were achiral and the stereochemistry was known. This choice was based on the fact that racemates are easily obtained, and the Supplementary information about the eutomer could be accessed after the enantiomeric purification of the active racemates. No specific processing of tautomers was performed (different tautomers of the same molecule have different fingerprints and likely different scores).

Model selection was performed using k-fold (k = 4) cross-validation. It concerned two parameters: the penalty parameter and the operating threshold probability. The penalty parameter was selected

TABLE 1 Statistical outlook of the initial dataset (each column represents an assay and the concentration at which compounds were tested)

Objectives	Activity	5-HT2A	5-HT2B	α1	D1	Nav 1.2	hERG	RLM	HLM	Caco-2 Fabs	Caco-2 efflux
Concentration	30 nM	10 μΜ	10 μΜ	10 μΜ	10 μΜ	10 μΜ	10 μΜ	-	-	-	-
Filled % ^a	29%	28%	26%	33%	28%	30%	59%	90%	90%	87%	77%
Blueprint threshold ^b	≥30%	≤50%	≤50%	≤50%	≤50%	≤50%	≤30%	≥50%	≥50%	≥90%	≤15
In blueprint rate ^c	59%	29%	35%	33%	53%	68%	45%	49%	35%	61%	80%

a "Filled %" describes the % of molecules in the dataset which have data in the assay.

to maximize the ROC AUC.³³ Once the penalty parameter selected, the operating threshold probability to predict conformity to the TPP (noted as 1 in Figure S1) was selected on the former k-folds to maximize precision to the detriment of recall, in order to reduce the risk of false positives. The best model, trained on 80% of the data (i.e., training set) was subsequently tested on the remaining 20% of the initial dataset (i.e., test set).

Classification models were selected rather than regression models due to their higher performances (results not shown).

2.4 | Generative model

As explained above, many architectures of molecular deep generative models have been reported to date. At the time this study was conducted (it was initiated in 2017), fewer architectures had been published. Molecule generation and property optimization strategies were inspired by Segler et al. which uses a deep RNN generator.¹³

2.4.1 | Molecule generation strategy

A deep RNN, and more precisely a deep LSTM of three hidden layers of size 512, was used to generate molecules represented as SMILES.^{12,34} The LSTM was first trained on the ChEMBL database, using teacher forcing,³³ to build a character-based language model for generating SMILES strings.¹³

It is reminded that the role of a language model p is to model the next character probability distribution given the sequence of previous characters:

$$p(x_{t+1}|x_1x_2...x_t) = LSTM(x_{t+1}|x_1x_2...x_t)$$

SMILES are generated by iteratively sampling the next character from its inferred past conditioned distribution $p(x_{t+1}|x_1x_2...x_t)$. Generating a SMILES starts and ends, respectively, with the special tokens of the vocabulary "START" and "END."

The SMILES in the ChEMBL database were transformed into their canonical achiral RDKIT version. No data augmentation by enumerating the different ways of writing a SMILES, nor by enumerating the

tautomeric forms of the same compound was performed. Thus trained, the LSTM language model generates achiral SMILES. Identical compounds can be generated with different writings of their SMILES. Tautomers of the same compound are generated as distinct molecules. Scheme 1 represents the architecture of molecule generation.

2.4.2 | Project dataset distribution learning

The LSTM trained on ChEMBL database has learnt to generate molecules belonging to ChEMBL chemical space. In order to be scored, generated molecules should stay near the applicability domain of the QSAR models. This applicability domain can be approximated by the structural similarity to the molecules of the initial dataset. Thus, the previous LSTM model was re-trained in teacher forcing on the project dataset.³⁵ This second training allows to zoom in the chemical space studied so that QSAR models can be applied.

2.4.3 | Molecule optimization strategy

The molecule optimization strategy that was used is named "Hillclimb-MLE." ^{13,23} It is an iterative process where the LSTM generative model is fine tuned in teacher forcing on an optimal set of SMILES that evolves over time as follows: at each step, this set of SMILES is updated by retaining the top 10% of compounds generated at the previous step (Figure 1).

The optimality ranking was established using a scalar reward function that combines 13 targets:

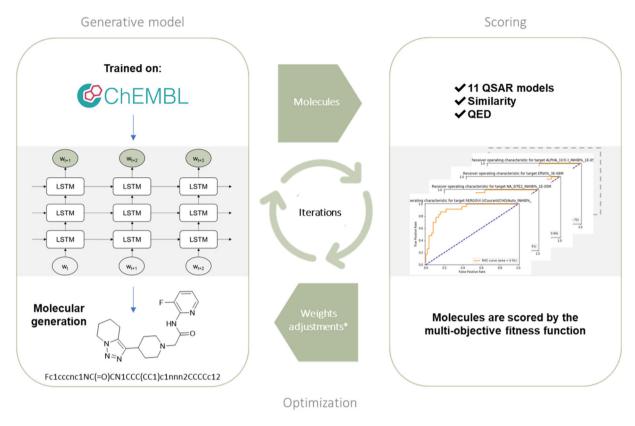
- + Eleven probabilities of activity $(p_i)_{1 \le i \le 11}$ returned by the classifiers described above (QSAR models built on the training data set):
 - + Similarity to the project dataset D, computed as:

$$\begin{split} S(\textit{mol}) &= \mathsf{max}\big(\big\{\mathsf{Tanimoto_Similarity}\big(\textit{mol}, \textit{mol}_j\big); \textit{mol}_j \in \mathsf{D}\big\}\big) \\ &+ \mathsf{QED}.^{36} \end{split}$$

Denoting $(x_i)_{1 \le i \le 13}$ and $(T_i)_{1 \le i \le 13}$, respectively, the vector of our 13 targets of interest and their thresholds for being in the blueprint, the reward function used in this project was the following:

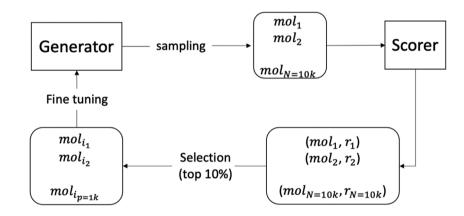
^bBlueprint threshold is the value set as the objective to achieve in each assay.

^cIn blueprint rate is the percentage of molecules meeting each objective individually.



SCHEME 1 Generative model architecture

FIGURE 1 Hill climbing procedure for optimizing the generator



$$Reward((x_i)_{1 \le i \le 13}) = -\sum_{i=1}^{13} \log \left(\frac{x_i}{T_i}\right)$$

Thresholds of QSAR scores (i.e., $(T_i)_{1 \le i \le 11}$) are their corresponding probability operating thresholds. The selected thresholds for similarity to the project dataset (i.e., T_{12}) and QED (i.e., T_{13}) are 0.5 and 0.4, respectively.

We could have relied on similarity as part of the reward function to learn the distribution of the initial dataset without re-training in teacher forcing. However, the generator would have needed more time to learn the initial distribution and the generated compounds in early steps of optimization would have been assigned predictions out of the applicability domain of the models.

Results are seed dependent and thus, many runs were conducted (10 runs), each one leading to new propositions to solve the problem.

2.5 | Assessment and ranking of generated compounds

Virtual candidates were ranked on their overall probability of being in the TPP, their Quantitative Estimate of Druglikeness (QED),³⁶ and their similarity to the initial dataset (i.e., Tanimoto distance). The applicability domain of the QSAR models is a critical point and must be carefully monitored to avoid false positives.

To help chemists assess the novelty and risk associated with the proposed molecules, a specific visualization was developed, by adapting the similarity map visualization.³⁷ This visualization, which we have named "applicability map" (Figure S2, Supplementary), enables to highlight, for each proposed molecule, the atoms which are either present or absent in the initial dataset, as follows: (a) in green, the atoms which are very well known because they appear very often in the same chemical environment in the initial dataset (i.e., the lead scaffold for instance); (b) in red, the new atoms or atoms already known but appearing in a new position; and (c) not highlighted: the atoms which have been seen before in the same position, but only a few times.

2.6 | Compound selection

From the newly generated library, the designed molecules were selected for synthesis and test based on their algorithmic ranking, structural novelty, synthetic accessibility, and consistency of the ADME predictions with those provided by global predictive models available at Servier.

3 | RESULTS

3.1 | Initial dataset analysis

The initial dataset, containing 881 molecules evaluated for 11 objectives, was sparse, with 10–70% missing data rates depending on the objectives. Due to the specificity of the primary assay, a complex ex vivo phenotypic assay, the ADME assays were very well documented, whereas only 251 compounds had been measured in the primary activity and selectivity assays. The dataset was well balanced, with >50% compounds meeting individually most objectives, with lower rates (\sim 30–35%) observed for 5-HT2A, 5-HT2B, alpha 1, and HLM.

The evolution of the percentage of compounds meeting each objective during the chronology of the project is displayed in Figure 2. It shows that the project team had been able to substantially increase the performance across iterations for Na_v1.2, hERG, RLM, HLM, with 80–90% of designed molecules meeting the required goal at the end of the program. Conversely, performance had strikingly dropped on 5-HT2A, alpha 1, D1, and permeability assays. As an example, only 6% of the last 50 molecules synthesized met the 5-HT2A selectivity objective. The colors in Figures 1 plot give an idea of the timeline of the project: (1) light gray for values found in molecules evaluated in the beginning of the project, molecules 1–780; (2) medium gray, for molecules developed based on the SAR for the preliminary results, molecule 781–830; (3) dark gray for late-stage molecules, expected to have the best profile based on the knowledge of more than 800 synthesized molecules, molecules 831–881.

In the subset of 48 molecules out of 881 which had been measured against all 11 objectives, the average number of objectives met was 6.4 out of 11. Among these, 6 molecules appeared to have a promising profile, meeting 9 objectives out of 11 (Figure 3). Molecule

FIGURE 3 Structure and biological profiles of the most promising lead molecule in the initial project dataset. The values in green correspond to the molecules active with the optimal threshold, the values in yellow correspond to the molecules active with the tolerated threshold; and the values in red correspond to the inactive molecules

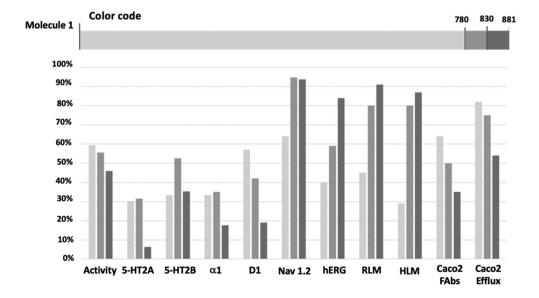


FIGURE 2 % of molecules in the initial project dataset meeting the different objectives along the chronology of the project (light gray: molecule 1–780; medium gray: molecule 781–830; dark gray: molecule 831–881)

TABLE 2 Biological profiles of the most promising lead molecules in the initial project dataset

Molecule ID	Activity	5-HT2A	5-HT2B	α1	D1	Nav 1.2	hERG	RLM	HLM	Caco-2 FAbs	Caco-2 Efflux
mol 732	194.0	20.0	18.0	1.0	4.0	0.0	19.0	82.85	63.35	88.99	26.2
mol 663	83.0	69.0	-25.0	45.0	6.0	13.0	6.4	69.04	31.93	97.6	1.96
mol 559	46.0	46.0	69.0	14.0	14.0	-14.0	25.8	60.28	25.43	98.86	0.75
mol 555	48.0	71.0	48.0	12.0	14.0	39.0	25.0	68.83	33.58	99.37	0.39
mol 550	115.0	76.0	15.0	37.0	-3.0	-13.0	5.4	80.82	83.54	72.24	12.3
mol 435	46.0	6.0	44.0	29.0	-11.0	20.0	12.4	93.11	78.36	73.8	34.1

Note: The colors correspond to the range of activity of the molecules. The values in green correspond to the molecules active with the optimal threshold, the values in vellow correspond to the molecules active with the tolerated threshold; and the values in red correspond to the inactive molecules.

732 (mol 732) was the best compound in the whole dataset, meeting all objectives except absorption, which was nearly met, and efflux.

It is worth noting that the 1,2-benzisoxazole in **mol 732** was also found in 61% of the project's compounds, and in 78% of the last 50 compounds made by the project team, indicating the importance that had been given by the medicinal chemistry team to that substructure, as a seemingly promising avenue for achieving a good balance between all desired properties. Only a couple of piperidine and piperazine linkers were used throughout the project, while more variability had been introduced in the East part heterocycles.

Also worthy of note, as shown in Table 2, a compound with a promising profile, mol 435, meeting 9 out of 11 objectives but missing absorption and efflux, quite close to mol 732 in terms of biological profile, had been obtained much earlier in the project. Two hundred and ninety-seven additional molecules were needed to partially improve the overall compound profile. During the design process from mol 435 to mol 732, permeability objectives were met in three molecules (mol 555, mol 559, and mol 663), but only to the detriment of 5-HT2A/B selectivity or metabolic stability.

3.2 | QSAR models

On average, the QSAR predictive models performed well with high precision in the test sets, except for 5-HT2B (precision 67%). Interpretability of the results was difficult for activity, alpha 1 and 5-HT2A due to the small number of positive compounds in the test set (confusion matrices are provided in Figure S1, ROC AUC plots are provided in Figure S3). The selected models were then trained on the whole dataset before switching to the generative phase of our work.

3.3 | Al-designed molecules

The generative algorithm designed 150 virtual compounds predicted to be optimal with regards to the project's TPP (i.e., predicted to meet

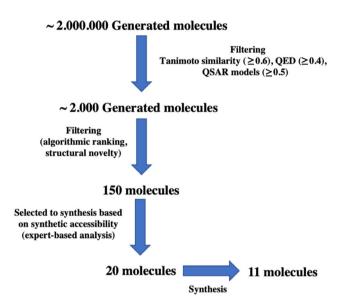


FIGURE 4 Pipeline of filtering AI generated molecules, from generation to synthesis

the required threshold for all targets), and with reasonable complexity as assessed by a chemist (at the time of the study, no satisfying synthetic accessibility scoring tool was available to help prioritize compounds). Among the 150 generated molecules, 20 were selected. In the 3-week timeframe allowed for compound synthesis: 11 compounds were successfully synthesized and tested on all the project's assays, whereas 9 molecules failed to be synthesized. Figure 4 shows an overview of the number of molecules filtered in each step, from generation to synthesis, and the molecules which were synthesized are represented in Figures 5 and 6.

After synthesis and test, the Al-generated candidates were found to outperform the initial library, including the last 50 compounds made within the project. The average number of objectives met by the Al-designed compounds was 9.5 (i.e., 86% success rate) versus 6.4 (i.e., 58% success rate) previously. Moreover, the Al-generated

096987x, 2022,

FIGURE 5 MPO profile of the best Al-designed molecule, **mol 885.** The values in green correspond to the molecules active with the optimal threshold

Caco-2 Efflux: 7.0

molecules reversed the decreasing trend in TPP conformity observed in the last molecules of the library (Figure 7A). Analysis shown in Figure 7B illustrates that, compared to the initial dataset, novel molecules were better on activity (i.e., in the blueprint 65% of the time) and excellent for all selectivity and permeability criteria (i.e., over 90% of the time in the blueprint). Metabolic stability, however, was lower, with a 55% conformity rate. More importantly, from the 11 new compounds, one met simultaneously all 11 objectives of the TPP (Figure 5) and two compounds met 10/11 objectives (Figure 6), while being just below the required threshold, within the error margin of the assay, for the missed objective.

The best Al-designed compound (mol 885), meeting all objectives, is represented in Figure 5. Notably, this compound contains a [1,2,3] triazolo[1,5-a]piperidine moiety which was very rare in the initial data set, appearing in only six molecules, and always correlated to poor permeability and efflux, which had led the project team to stop investigating this motif. It is remarkable that the Al algorithm retained that substructure, combining it with a 3-fluoropyridine in the East part, which had never been tried before. Surprisingly, the association of this discarded substructure with an unexplored heterocycle turned out to be a winning combination for solving the MPO objective of the project.

As a matter fact, the 11 Al-designed compounds that were synthesized and tested displayed functional groups that were either rare in the initial dataset or never tried earlier in the project (see Figure 6). It suggests that this method can propose significant innovations, by its ability to identify favorable modifications, even with few data to learn from.

One striking example is **mol 886**, where an aliphatic group was introduced in replacement of an aryl moiety, where only aromatic moieties had been used before at this specific position.

The AI algorithm was also able to optimize ADME properties in sub-series with specific issues. For example, it was able to design permeable compounds within the 6,7-dihydro-4H-triazolo[5,1-c][1,4] oxazine sub-series while maintaining safety and stability, when all compounds in that sub-series had permeability issues. Likewise, within the pyrido-isoxazole series, compounds with reduced efflux were identified while maintaining safety and stability (Figure 8).

An analysis of the drug-likeness profile of the compounds based on their property forecast index (PFI), molecular weight (MW) and sp3 fraction was performed.³⁹ The plot of PFI vs MW is presented in Figure 9. Ten out of 11 Al-designed compounds were found to have a very favorable profile with low PFI, low MW, and high sp3 fraction, compared to the molecules from the initial data set.

To provide insights about structural diversity and chemical space features of both the initial dataset and Al designed compounds, a principal component analysis (PCA) was computed on the Morgan fingerprints (i.e., extended connectivity fingerprints [ECFP] of 1024 bits, radius 2) of the molecules in the dataset.³⁸ First, a representation of the 251 compounds from the initial library that were measured in the primary activity assay is provided (Figure 10). This plot reveals the absence of a probability gradient or narrow area of activity since active molecules can be found in all areas of the explored chemical space. Conversely, a display of the number of objectives met by these 251 molecules allows to delineate an area where the MPO score is the highest (i.e., the upper left corner of the plot in Figure 9).

Strikingly (Figure 10), the AI algorithm did not design any molecule in that seemingly promising chemical space. All AI-designed structures are indeed located in a distinct yet specific area, demonstrating the capacity of this algorithm to come up with non-trivial solutions.

4 | DISCUSSION

A typical hurdle of MPO is that optimization of some objectives leads to a drop of performance in others, but the present method allowed the design of compounds that were simultaneously optimized on 11 parameters.

Yet, several features of the initial dataset were key to enable achieving such performance.

Overall, the performances of the models built to predict bioactivity on each assay were good, thereby validating the approach of project data-guided optimization. This requires enough data to build a decent model (in our case, the least documented assay had $\sim\!250$ data points) and a reasonably well-balanced data set with enough compounds meeting each objective individually. Also, the generative model was able to find theoretical solutions to the MPO challenge within the chemical space of the project, meaning that based on the available data, there were indeed ways to solve the apparent anticorrelations between the objectives.

This favorable configuration may not be present in all cases, and the potential of the method to solve MPO challenges in more complex cases remains to be demonstrated. Several approaches could be envisaged to circumvent the lack of balanced data on some objectives, such as using generic models trained on large and diverse legacy data, that is, for ADME properties prediction, or using structure-based modeling to guide optimization on target or anti-targets if such structural information is available. To address the tricky issue of the inability to identify structures solving the MPO challenge within the project's chemical space, adding an active learning component to the

Anti-targets: 6/6

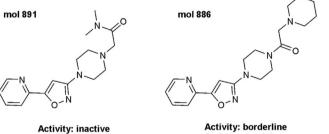
ADME: 3/4

FIGURE 6 Structures and biological features of original compounds sampled by the DL algorithm. Borderline compounds have been measured below the desired activity threshold but within the error margin of the assay while active and inactive molecules were measured, respectively, above and below the threshold

Presence of a 4,7-dihydro-5H-pyrano[4,3-d]pyrazole moiety which appears 5 times in the initial dataset

Anti-targets: 6/6

ADME: 1/4



Activity: inactive Anti-targets: 6/6 ADME: 4/4

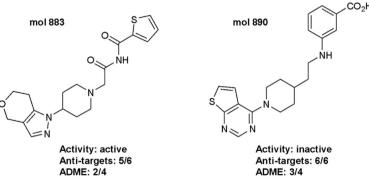
Anti-targets: 6/6

ADME: 2/4

Activity: borderline Anti-targets: 6/6 ADME: 4/4

Presence of a 2-(1,2-oxazol-5-yl)pyridine moiety which appears 13 times in the initial dataset First introduction of an aliphatic group in the East part of the molecules

Presence of a 6,7-dihydro-4H-triazolo[5,1-c][1,4]oxazine moiety which appears 4 times in the initial dataset



Presence of a 6,7-dihydro-4H-pyrano[3,4-d]pyrazole moiety which was present 5 times in the initial datasset

Presence of a thieno[2,3-d]pyrimidine moiety which was absent from the initial datasset

SAR models to guide the optimization and/or using more generic molecular representations (2D or 3D pharmacophore-based) to build the QSAR predictive models could be considered.

All in all, it is still uncertain to which extent Al-guided optimization can bring benefit to Lead Optimization in terms of reduction of number of compounds and number of iterations needed to identify a

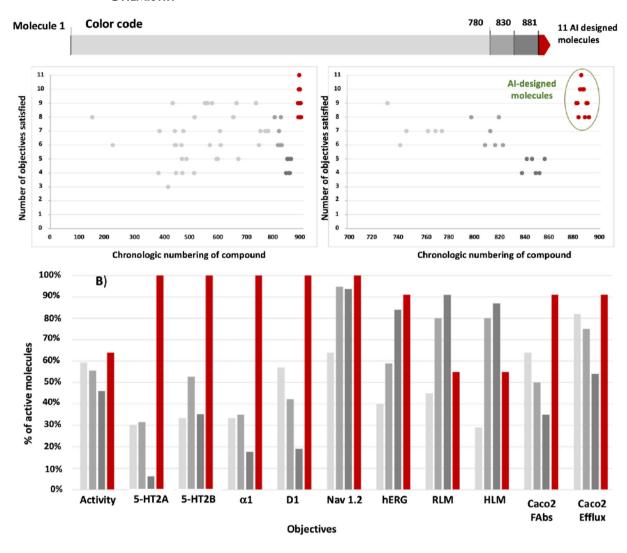


FIGURE 7 (A) Number of objectives satisfied according to project's chronology (vertical axis: number of objectives satisfied /horizontal axis: chronologic numbering of compound/please note initial data was sparse with only 48 compounds tested on all criteria. (B) Hit rate comparison between Al-designed candidates and initial molecules for each TPP objective. (light gray: molecule 1–780; medium gray: Molecule 781 to 830; dark gray: molecule 831 to 881; red: 11 Al-designed molecules)

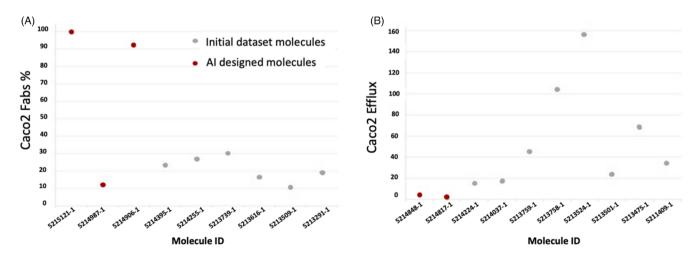


FIGURE 8 Permeability (A) and efflux (B) properties of Al-designed vs original dataset compounds in the (A) 6,7-dihydro-4H-triazolo[5,1-c] [1,4]oxazine series (left) and (B) pyrido-isoxazole series (right)

1096987x, 2022, 10, Downloaded from https://onlinelibrary.wiley.com/doi/10.1002/jcc.26826 by Cochrane France, Wiley Online Library on [15/05/2025]. See the Terms

on Wiley Online Library for rules of use; OA articles are governed by the applicable Creative Commons License



FIGURE 9 Plot of MW in function of PFI for initial molecules and Al-designed compounds

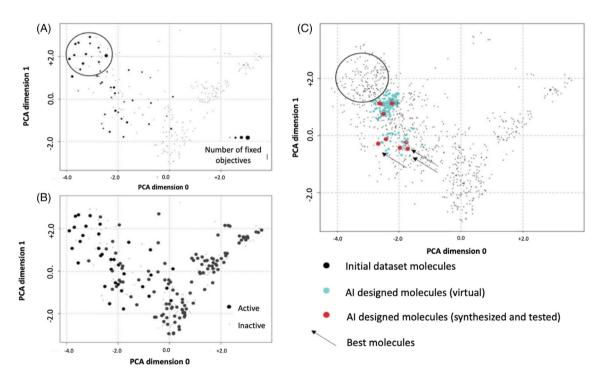


FIGURE 10 (A) PCA of the same 251 active compounds with correlation to the TPP criteria hit rates. (B) PCA of the 251 active compounds from the initial dataset. (C) Plot of the Al-designed molecules

new molecular entity, and it would be of high interest to test such approach along several design-make-test cycles, starting early in the Lead Optimization phase, to assess the magnitude of the benefit brought by Al. Ideally, this should be conducted as a comparative "blind" study comparing the Al approach to a traditional approach to enable to draw strong conclusions.

Also, worthwhile mentioning, the selection process of the Aldesigned molecules was not only based on pure data-driven ranking. Molecules were selected based on their scores on the predictive models, but also based on their synthetic accessibility and the expert input of medicinal chemists and computational chemists using their expertise as well as specific data visualization tools to remove poor

096987x, 2022,

quality compounds or potential false positives. This selection process associating human expertise and data visualization to rank and select Al-driven ideas was probably an important success factor in this project. Indeed, although not addressed in this paper, issues with synthetic accessibility, complexity, structural alerts issues, or sheer meaninglessness of certain Al propositions did occur in this project, although they remained minor in this context. These issues currently prevent a fully automated compound selection and rather advocate for a collaboration between chemists and Al, enabling to get the best of both worlds. Recently, notable progress has been made in the development of efficient methods for high throughput synthetic accessibility scoring³⁹⁻⁴² which opens the perspective of an increased automation of the process.

Besides accelerating the discovery of active molecules with a good MPO profile, another value of the approach was to open up new chemical space, in a phase when the project team apparently had already "homed in" on a fairly well-defined scaffold. The association of the [1,2,3]triazolo[1,5-a]piperidine moiety, very rare in the initial data set and correlated to poor permeability and efflux, with a 3-fluoropyridine, never tried before, was the answer to solve the MPO problem.

5 | CONCLUSION

Exploiting a sparse dataset of 881 molecules measured on 11 bioactivity assays, a DL-based AI de novo design algorithm was able to generate 150 virtual compounds with optimal in silico profiles against all desired characteristics of the project's TPP. Among those, 11 compounds were synthesized and measured on all 11 criteria of the TPP. The AI-designed molecules outperformed the ones designed by traditional medicinal chemistry approaches, achieving superior MPO scores. More importantly, three of those were found to meet the project's TPP, one of them strictly meeting all MPO objectives, the other two matching 10 objectives and being in the error margin of the assay for the last one. The AI algorithm came up with functional groups, which, although being rare or absent in the initial dataset, turned out to be highly beneficial for the MPO.

To our knowledge, this is the first report of a successful application of deep learning to de novo design for solving an MPO issue in an actual drug discovery project, moreover on a large number of objectives. This brings unequivocal evidence of the potential of this technology to bring substantial improvements to medicinal chemistry. The use of such an approach in earlier stages of drug discovery (i.e., hit discovery, hit to lead and early LO) is under investigation. Improvement needs have been identified and are being addressed, notably regarding synthetic accessibility, compound complexity and domain of applicability of the predictive models.

ACKNOWLEDGMENTS

The authors would like to thank Dr. Vinicius Barros Ribeiro da Silva for his help on the redaction and revision of the manuscript.

DATA AVAILABILITY STATEMENT

Research data are not completely presented due to the commercial interest in the non-patented molecules. However, a sample of the best molecules from the initial dataset and the Al-generated molecules synthesized and tested are shared.

ORCID

Yann Gaston-Mathé https://orcid.org/0000-0002-9635-7324

REFERENCES

- [1] R. S. Bohacek, C. McMartin, W. C. Guida, Med. Res. Rev. 1996, 16, 3.
- [2] C. Lipinski, A. Hopkins, Nature 2004, 432, 855.
- [3] C. A. Nicolaou, N. Brown, Drug Discov. Today Technol. 2013, 10, e427
- [4] N. C. Firth, B. Atrash, N. Brown, J. Blagg, J. Chem. Inf. Model. 2015, 55, 1169.
- [5] S. M. Paul, D. S. Mytelka, C. T. Dunwiddie, C. C. Persinger, B. H. Munos, S. R. Lindborg, A. L. Schacht, *Nat. Rev. Drug Discov.* 2010, 9, 203.
- [6] A. Talevi, Methods Mol. Biol. 2018, 1762, 1.
- [7] D. Kitchen, H. Decornez, J. Furr, J. Bajorath, Nat. Rev. Drug. Discov. 2004, 3, 935.
- [8] A. Lavecchia, C. Di Giovanni, Curr. Med. Chem. 2013, 20, 2839.
- [9] B. J. Neves, R. C. Braga, C. C. Melo-Filho, J. T. Moreira-Filho, E. N. Muratov, C. H. Andrade, Front. Pharmacol. 2018, 9(1275), 1.
- [10] C. A. Nicolaou, J. Apostolakis, C. S. Pattichis, J. Chem. Inf. Model. 2009, 49, 295.
- [11] P. Schneider, W. P. Walters, A. T. Plowright, N. Sieroka, J. Listgarten, R. A. Goodnow, J. Fisher, J. Jansen, J. Duca, T. Rush, M. Zentgraf, J. E. Hill, E. Krutoholow, M. Kohler, J. Blaney, K. Funatsu, C. Luebkemann, G. Schneider, *Nat. Rev. Drug. Discov.* 2020, 19, 353.
- [12] D. Weininger, J. Chem. Inf. Comput. Sci. 1988, 28, 31.
- [13] M. H. S. Segler, T. Kogej, C. Tyrchan, M. P. Waller, ACS Cent. Sci. 2018, 4, 120.
- [14] M. Krenn, F. Häse, A. Nigam, P. Friederich, A. Aspuru-Guzik, Mach. Learn. Sci. Technol. 2019, 1(45024), 1.
- [15] A. Nigam, P. Friederich, M. Krenn, A. Aspuru-Guzik, arXiv, 2019. Doi: arXiv:1909.11655v4.
- [16] N. O'Boyle, A. Dalke, ChemRxiv, 2018. https://doi.org/10.26434/ chemrxiv.7097960.v1.
- [17] Z. Zhou, S. Kearnes, L. Li, R. N. Zare, P. Riley, ArXiv, 2019. Doi: arXiv: 1810.08678v3.
- [18] M. Popova, M. Shvets, J. Oliva, O. Isayev, ArXiv, 2019. Doi: arXiv: 1905.13372.
- [19] W. Jin, R. Barzilay, T. Jaakkola, ArXiv, 2018. Doi: arXiv: 1802.04364v4.
- [20] R. Gómez-Bombarelli, J. N. Wei, D. Duvenaud, J. M. Hernández-Lobato, B. Sánchez-Lengeling, D. Sheberla, A. Aspuru-Guzik, ACS Cent. Sci. 2018, 4, 268.
- [21] G. Guimaraes, B. Sanchez-Lengeling, C. Outeiral, P. Cunha Farias, A. Aspuru-Guzik, ArXiv, 2017. Doi: arXiv:1705.10843v3.
- [22] M. Olivecrona, T. Blaschke, O. Engkvist, H. Chen, Aust. J. Chem. 2017, 9(48), 1.
- [23] D. Neil, M. H. Segler, L. Guasch, M. Ahmed, D. Plumbley, M. Sellwood, N. Brown, ICLR, 2018.
- [24] M. J. Kusner, B. Paige, J. Hernández-Lobato, ArXiv, 2017. Doi: arXiv: 1703.01925v1.
- [25] R. Winter, F. Montanari, A. Steffen, H. Briem, F. Noé, D.-A. Clevert, Chem. Sci. 2019, 10, 8016.
- [26] G. Schneider, Nat. Rev. Drug Discov. 2018, 17, 97.
- [27] K. Gao, D. D. Ngugyen, M. Tu, G. W. Wei, J. Chem. Inf. Model. 2020, 60, 5682.

- [28] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, E. Duchesnay, J. Mach. Learn. Res. 2011, 12, 2825.
- [29] J. Bergstra, D. Yamins, D. D. Cox, Proc. of the 30th Int. Conf. Machine Learn. (ICML 2013) 2013, 28, 115.
- [30] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, R. Jozefowicz, Y. Jia, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, M. Schuster, R. Monga, S. Moore, D. Murray, C. Olah, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, X. Zheng, ArXiv, 2016. Doi: arXiv: 1603.04467.
- [31] G. Landrum, RDKit: Open-source cheminformatics, 2006.
- [32] D. Rogers, M. Hahn, J. Chem. Inf. Model. 2010, 50, 742.
- [33] A. P. Bradley, Pattern Recognit. 1997, 30, 1145.
- [34] S. Hochreiter, J. Schmidhuber, Neural Comput. 1997, 9, 1735.
- [35] R. J. Williams, D. Zipser, Neural Comput. 1989, 1, 270.
- [36] G. R. Bickerton, G. V. Paolini, J. Besnard, S. Muresan, L. Andrew, A. L. Hopkins, Nat. Chem. 2012, 4, 90.
- [37] H. Eckert, J. Bajorath, Drug Discovery Today 2007, 12, 225.

- [38] S. Wold, K. Esbensen, P. Geladi, Chemom. Intell. Lab. Syst. 1987, 2, 37.
- [39] C. W. Coley, L. Rogers, W. H. Green, K. F. Jensen, J. Chem. Inf. Model. 2018, 58, 252.
- [40] A. Thakkar, V. Chadimova, E. J. Bjerrum, O. Engkvist, J. L. Reymond, Chem. Sci. 2021, 12, 3339.
- [41] H. Tajmouati, M. Parrot, A. Skiredj, R. Fourcade, N. D. Huu, Q. Perron, Y. Gaston-Mathé, Royal Chem. Soc. Artif. Intellig. 2020.
- [42] https://spaya.ai

SUPPORTING INFORMATION

Additional supporting information may be found in the online version of the article at the publisher's website.

How to cite this article: Q. Perron, O. Mirguet, H. Tajmouati,

- A. Skiredj, A. Rojas, A. Gohier, P. Ducrot, M.-P. Bourguignon,
- P. Sansilvestri-Morel, N. Do Huu, F. Gellibert,
- Y. Gaston-Mathé, J. Comput. Chem. **2022**, 43(10), 692. https://doi.org/10.1002/jcc.26826