



# Active-learning strategies in computer-assisted drug discovery

Daniel Reker and Gisbert Schneider

Q1

Swiss Federal Institute of Technology (ETH), Department of Chemistry and Applied Biosciences, Vladimir-Prelog-Weg 4, 8093 Zürich, Switzerland

High-throughput compound screening is time and resource consuming, and considerable effort is invested into screening compound libraries, profiling, and selecting the most promising candidates for further testing. Active-learning methods assist the selection process by focusing on areas of chemical space that have the greatest chance of success while considering structural novelty. The core feature of these algorithms is their ability to adapt the structure–activity landscapes through feedback. Instead of full-deck screening, only focused subsets of compounds are tested, and the experimental readout is used to refine molecule selection for subsequent screening cycles. Once implemented, these techniques have the potential to reduce costs and save precious materials. Here, we provide a comprehensive overview of the various computational active-learning approaches and outline their potential for drug discovery.

## Introduction

Q2

The concept of iterative molecular design, synthesis, and testing forms a central pillar of drug discovery; it provides the basis for our understanding of the underlying structure–activity relation (SAR). Iterative synthesize-and-test cycles with SAR model adaptation to newly obtained activity data improve the overall quality of the designer compounds and help reduce experimentation costs. Similarly, the screening of existing compounds profits from such feedback-driven picking: within a fixed budget, adaptive screening rounds through multiple acquisition-and-test cycles can lead to significantly better solutions compared with a single large screen [1,2]. The crucial step in each learning cycle is the formulation of a well-motivated hypothesis for compound generation (*de novo* design) or compound picking (when screening from a compound pool) based on the available SAR data. The selected molecules can either be hypothesized actives or readily available compounds that will improve the model by elucidating poorly understood parts of the SAR. Commonly, an interdisciplinary team of scientists generates the new hypothesis by inferring from their expertise and medicinal chemistry ‘intuition’. Therefore, any design hypothesis is easily biased towards preferred chemistry [3,4] or pre-disposed model interpretation [5,6]. Although expert knowledge is

indisputably important for successfully guiding drug discovery projects, an unbiased perspective during the compound selection process can lead to structurally surprising chemical agents with the desired novelty, bioactivity, and physicochemical properties [7]. Moreover, with the recent advances of microfluidics-assisted integrated medicinal chemistry platforms (e.g., lab-on-a-chip systems [8]), the generation of an accurate and suitable molecular design hypothesis and, consequently, the selection of new compounds for synthesis and testing, becomes the bottleneck in an otherwise automatable optimization process [9].

Computational models act as rapid and objective decision makers in this decisive selection step (Fig. 1a) [10,11]. Active learning (also known as ‘selective sampling’) is an umbrella term from the field of machine learning for methods that select data points for testing and feeding back into the model [12,13]. Approximately 15 years ago, the term was introduced to drug discovery [14]. Recently, the topic has gained momentum, driven by technological advancements in small-scale organic synthesis systems and the accuracy of machine-learning prediction models. Here, we provide a comprehensive overview of investigations that have applied active-learning techniques to drug discovery. We focus on methods for finding novel chemical structures and discuss possible future directions of algorithm development and how these might help solve current challenges in computer-assisted drug design (Box 1).

Q3

Corresponding author: Schneider, G. (gisbert.schneider@pharma.ethz.ch)

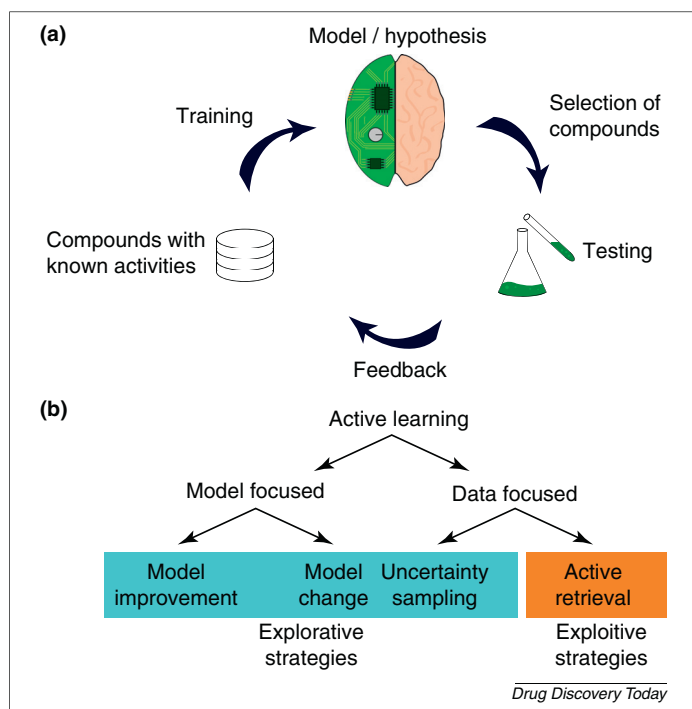


FIG. 1

(a) Schematic of the active-learning concept. Known activity data are provided as training data to a machine-learning model that generalizes this knowledge. A selection strategy is used that picks from a list of new molecules with unknown activity. These selection strategies usually try to identify molecules that would be particularly suited for improving the model quality (explorative strategies) if they are included in the training database with their activity value. Alternatively, molecules are selected that might have favorable activity values (exploitive strategies). After the selected molecules have been tested ('labeled'), they are added to the training data to train an improved machine-learning model. (b) Conceptual comparison of different active-learning strategies. These can be distinguished methodologically according to whether the selection strategy is derived from the whole model ('model focused') or by examining individual data points ('data focused'). When compounds are selected with the whole model in mind, the strategies are explorative. Possible implementations are predicting or calculating the change in model architecture ('model change') or the improvement of the model ('model improvement'; e.g., variance reduction or error on the test set). When examining individual data points, models can either be exploitive ('active retrieval') or use the error or uncertainty on the individual data points to perform confined model optimization ('uncertainty sampling').

## Exploration versus exploitation

Compound selection strategies can be distinguished according to their underlying motivation (Fig. 1b): whereas some algorithms utilize the available information to retrieve compounds with certain properties ('exploitation'), others seek to improve the model by adding knowledge ('exploration'). From a technical point of view, exploration can either be performed from a molecule-centric perspective ('uncertainty sampling', i.e., selecting molecules that are predicted with low confidence by the model) or by explicitly estimating the impact of adding the additional data point on the error or architecture of the model ('model-centric' approaches). Explorative strategies sample more diverse chemical structures and rapidly increase the knowledge for the model (Fig. 2a), while not always proposing favorable structures in terms of their activity (Fig. 2b). Conversely, exploitive strategies retrieve active compounds with a greater probability, but do not

### BOX 1

#### Pseudocode for performing a retrospective active-learning investigation ('ActiveLearning')

```
function ActiveLearning(M,s):
  T, L, E ← splitStratified(M)
  m ← trainRFModel(T)
  for I ← 1 to 100 do:
    selected_mol ← s(L,m)
    L ← L \ {selected_mol}
    T ← T ∪ {selected_mol}
    m ← trainRFModel(T)
    evaluate(m,E)
  end do

function random(L,m):
  return pickOneRandom(L)

function exploitive(L,m):
  predictions ← m.predictActivity(L)
  return L[argmax(predictions)]

function explorative(L,m):
  uncertainty ← m.predictionVariance(L)
  return L[argmax(uncertainty)]
```

The function takes descriptions and activities of a set of molecules ('M') and a selection function ('s') that is used for the picking of molecules. First, the molecular data are split into three subsets in a stratified manner according to activity. Afterwards, the training data ('T') are used for initial model training ('trainRFModel'). The active learning is performed for 100 iterations in which we first pick a molecule from the learning data ('L') according to the selection function. This selected molecule is then removed from the learning data and added to the training data, with which the model is retrained. The performance of the new model can then be evaluated ('evaluate'), for example according to the error on the test data ('E', Fig. 2a, main text), the activity of the picked molecule (Fig. 2b, main text), or the number of scaffolds known to the model (Fig. 2c, main text). As examples of selection functions, we show pseudocode for a random strategy ('random') that picks a random molecule from the set, an exploitive strategy ('exploitive') that picks the molecule with the highest predicted activity, and an explorative strategy ('explorative') that picks the molecule with the highest prediction uncertainty (e.g., the maximum variance according to the individual activity predictions of the trees of the random forest model).

necessarily add knowledge to the model. In fact, the model quality can even decrease over time when using an exploitive strategy because of the introduction of a strict bias towards highly active compounds (Fig. 2a). Various strategies for either of the two compound selection principles have been proposed and validated in the context of drug discovery (Table 1).

Explorative approaches have proven particularly attractive when aiming at novel chemotypes with desired bioactivities (Fig. 2c). For example, to probe for the applicability of uncertainty sampling to explorative drug design, Lemmen and coworkers developed both a jury of Perceptrons and a support vector machine (SVM) model to distinguish thrombin ligands from 'inactives' [14,15]. Model optimization was conducted by adding examples

Download English Version:

<https://daneshyari.com/en/article/10885969>

Download Persian Version:

<https://daneshyari.com/article/10885969>

[Daneshyari.com](https://daneshyari.com)