

Review

Molecular docking and pharmacophore modelling; a bridged explanation with emphasis on validation

Rahma SR. Mahrous ^{*a}, Hoda M. Fathy ^a, Rasha M. Abu EL-Khair ^a, Abdallah A. Omar ^a and Reham S. Ibrahim ^a

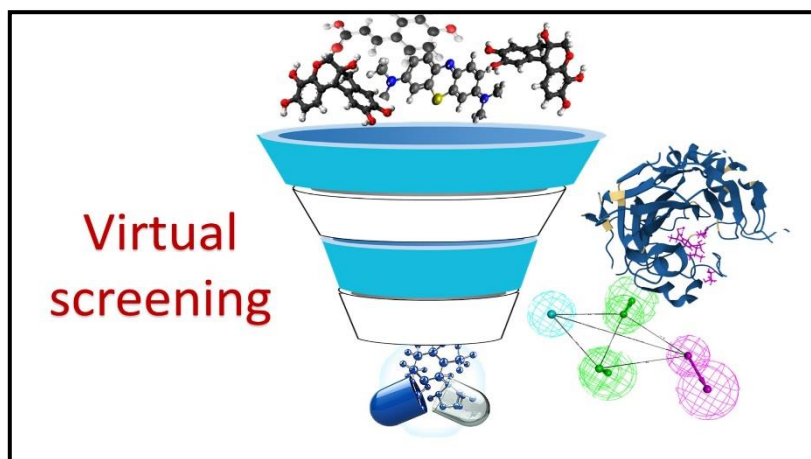
^a Department of Pharmacognosy, Faculty of Pharmacy, Alexandria University, Egypt

***Corresponding author.** Department of Pharmacognosy, Faculty of Pharmacy, Alexandria, Egypt. 1 El-Khartoum Square - Azarita, Postal code: 21521

E-mail: rahma.mohammed@alexu.edu.eg

Abstract:

Virtual screening (VS) techniques have emerged in the past decade as an efficient strategy in lead identification. Molecular docking as well as pharmacophore and 3D-QSAR modelling are two major VS techniques standing out as cornerstones in the process of new drug discovery. Explanation of the main virtual screening techniques as well as the



most commonly used validation parameters are discussed thoroughly emphasizing their use and shortcomings. Criteria for the selection of benchmarking datasets and training sets for molecular docking, pharmacophore, and 3D-QSAR models are discussed. Understanding the basics behind these techniques and their validation is crucial to judge the validity of the obtained results. Computational technologies have witnessed great improvement in the last few years which had a great impact on the improvement of virtual screening. Fields such as cloud computing and deep learning algorithms are among the technologies modeling the future of computer-aided drug design. This mini-review gives summarized knowledge for guiding beginners as well as experts in the field of virtual screening.

Keywords: Molecular docking, validation, 3D-QSAR, deep learning, cloud computing.

Received 7 December 2023

Accepted 30 January 2024

Published 30 January 2024

1. Introduction

In the past 100 years and with the development of molecular sciences, the discovery of new proteins and their unique structure has been attracting significant attention as a valuable source of lead drug targets ⁽¹⁾. Those leads can be synthetic, natural, or semi-synthetic compounds from a naturally occurring nucleus.

Computer-aided approaches (in-silico) are playing a crucial role in the process of lead drug identification complementing other techniques such as high throughput screening (HTS) bioactivity-guided fractionation, bioautography, HPLC- based activity profiling, MS & NMR based methods, and many others are examples of these approaches ^(2, 3). These computational methods are becoming increasingly popular due to their cost-effectiveness and time efficiency compared to conventional methods used in the drug discovery process. Computer-aided drug discovery (CADD) approaches incorporate different fields such as drug target identification, virtual screening of chemical structure databases for prototype identification, in silico assessment of pharmacokinetics properties and toxicity profile (ADMET) of drug candidates, and optimization of these candidate compounds ⁽⁴⁻⁶⁾.

Among CADD, virtual screening (VS) approaches have a significant impact on identifying lead compounds for a given target. Molecular docking, pharmacophore modelling, structure-activity relationship (SAR), and quantitative structure-activity relationship (QSAR) are utilized in virtual screening ^(7, 8). Being performed on a virtual level, validation of these techniques is the milestone in evaluating their outcomes. A massive achievement in computational technologies and the introduction of cloud computing and advances in machine and

deep learning led to the improvement in the field of virtual screening ⁽⁹⁾.

This mini-review is a simplified explanation of commonly used validation parameters and metrics that are applied in molecular docking and 3D-QSAR modelling. We also shed light on the basics that help researchers and even students to understand these techniques, their applications, and threats. In addition to some of the future perspectives of virtual screening techniques.

2. Virtual screening (VS) programs: a useful means for performing in-silico work

Generally, virtual screening (VS) programs are used for the following purposes: 1) identification of potential ligand–receptor interaction points, 2) reducing huge compounds database to the smaller matrix of predicted active compounds (hits), 3) designing novel compounds via using fragments and/or functional groups into new chemotypes, and 4) optimization of lead compound characters as pharmacokinetic properties and drug metabolism ⁽¹⁰⁾. Advances over the past 50 years resulted in the development of more than 60 different docking tools and programs such as Glide, GOLD 1.1, FlexX 1.8, DOCK, Autodock, and others ⁽¹¹⁾. A list of the commonly used docking softwares are given in **Table 1** ^(12, 13). In addition, a wide variety of pharmacophore modeling software like; Catalyst, Phase, HipHop, HypoGen, DISCO, MOE, and others were introduced, and are introduced in **Table 1** ^(14, 15). These softwares and engines differed mainly in the algorithms used for the search and filtering methods. These algorithms perform the general filters of the compound libraries. Identify the potential ligand–receptor interaction points. In addition, evaluate these interactions and determine their strength through scoring function ^(16, 17).

Table 1: Some molecular docking and pharmacophore modelling softwares of common use in virtual screening

Software	Reference	Website (if available)
Molecular Docking (SBVS)		
AUTODOCK	(18)	https://autodock.scripps.edu/
Dock	(19)	http://dock.compbio.ucsf.edu/
FlexX	(20)	http://www.biosolveit.de/flexx/
FRED	(21)	http://www.eyesopen.com/oedocking
GLIDE	(22)	http://www.schrodinger.com/
GOLD	(23)	https://www.ccdc.cam.ac.uk/solutions/software/gold/
ICM	(24)	http://www.molsoft.com/docking.html
Ligand-Receptor Docking	(12)	http://www.chemcomp.com/software-sbd.htm
Maestro	(25)	http://www.schrodinger.com/downloadcenter/
SLIDE	(26)	https://kuhnlab.natsci.msu.edu/software/slide/index.html
Surflex	(27)	https://www.biopharmics.com/
Virtual Docker	(28)	http://molexus.io/molegro-virtual-docker/
ZDOCK	(29)	http://zlab.umassmed.edu/zdock/
Pharmacophore modeling		
Catalyst	(30)	-
COMSIA	(31)	-
DISCO	(32)	-
GALAHAD	(33)	-
GASP	(34)	-
HipHop	(30)	-
HypoGen	(35)	-
LigandScout	(36)	https://www.inteligand.com/ligandscout/
MOE	(14)	https://www.chemcomp.com/Products.htm
PHASE	(37)	https://newsite.schrodinger.com/platform/products/phase/
PhDOCK	(38)	-
ZincPharmer	(39)	http://zincpharmer.csb.pitt.edu/

3. Molecular docking (Structure-based virtual screening)

Molecular docking uses computer technology to fit the conformations of ligands (usually small molecules or another protein) into protein binding sites and predict the best structure of receptor-ligand complexes that complements the protein binding site with the lowest energy⁽⁴⁰⁻⁴²⁾. The 3D structure of the target should be available for molecular docking and hence it is named structure-based virtual screening (SBVS). X-ray

crystallography and NMR studies help in solving the 3D structures,⁽⁴³⁾ or built using homology modelling with related proteins⁽⁴⁴⁾ and they are freely available at the Protein Data Bank (PDB) (<https://www.rcsb.org/>) facilitating molecular docking studies. In selecting this 3D structure, crystal structures with the bound ligands are preferred since they give information about features of the binding site and ligand-target interactions⁽⁷⁾. Two fundamental steps are performed in molecular docking software, first, prediction

of ligand orientation and conformation inside the target's active site (Pose prediction) ^(16, 45, 46). Second, the scoring or evaluation function measures the fitness of each predicted ligand pose into the target binding site through the calculation of their free binding energies where lower values correspond to more favorable ligand binding (ranking) ⁽⁴¹⁾.

4. Docking validation; method for evaluating docking performance

Validation of docking studies represents the evaluating tool to assess the accuracy of docking software ⁽⁴⁷⁻⁵¹⁾. Docking validation studies are divided into two interlacing classes, (i) Ligand pose prediction, and (ii) Ligand screening ^(41, 50) (**Fig. 1**).

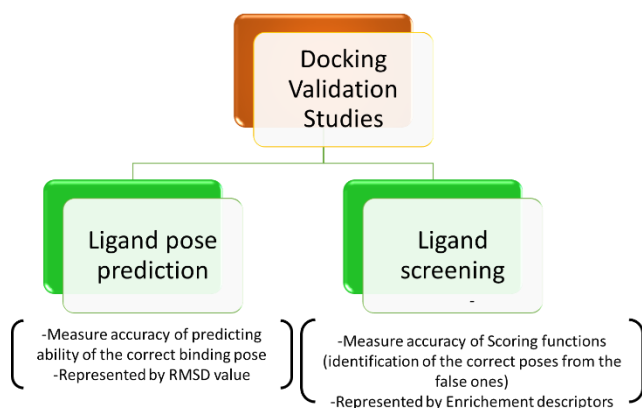


Fig. 1: Molecular docking validation studies

4.1. Ligand Pose Prediction

It is based upon the capability to replicate the experimental coordinates of previously docked ligands. Docking software packages are utilized to re-dock the original co-crystallized ligand into the target's active site. Then, computing the root mean square deviation (RMSD) that denotes the mean distance of the ligand's atoms in the model compared to those in the cocrystallized structure after receptor superposition (**Fig. 2**). Programs able to return poses with a preselected RMSD value (usually 1.5 or 2 Å) are accurate in pose prediction. It is accepted that RMSD values of 2 °A are an indication

of the accuracy of protein–ligand docking algorithms beyond which the prediction is considered not reliable ^(16, 49).

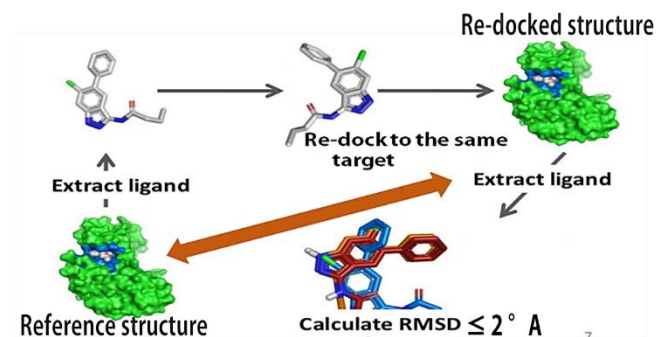


Fig. 2: Ligand pose prediction

4.2. Ligand Screening (Enrichment descriptors)

The key for molecular docking is the capability of docking protocol to pick active molecules among a set of inactives. This is performed through scoring and ranking that are evaluated by calculating enrichment descriptors. These descriptors indicate the capability of scoring functions to envisage the correct poses ^(16, 49). The VS technique is considered successful based on its ability to organize the active compounds in descending order of activity in a ranked hit list. The top fraction of the obtained hit list is likely to be evaluated using experimental *in-vitro* / *in-vivo* models ⁽⁵²⁾.

The two most important values in the majority of the enrichment descriptors are sensitivity and specificity. Sensitivity (Se, true positives) denotes the ratio between the number of active compounds recovered by the virtual screening method and the total number of active molecules in the database. Specificity (Sp, false positives), denotes the ratio between the number of inactive compounds that were not selected by the VS method and the total number of inactive present in the database ⁽⁵²⁾. These two terms are calculated according to the following equations*:

$$Se = \frac{TP}{TP+FN} \quad (1)$$

$$Sp = \frac{TN}{TN + FP} \quad (2)$$

*Where (TP): true positive compounds, (FN): false negatives that are active molecules not identified by the VS method. (TN): true negatives, the unselected inactive compounds, (FP): false positive compounds (decoys).

Evaluation of VS softwares is carried out using several metrics such as Area Under the Curve (AUC), Enrichment factor (EF), Boltzmann-enhanced discrimination of ROC (BEDROC), Robust initial enhancement (RIE), Receiver operating characteristic (ROC), and others⁽⁵³⁾. It is of great importance that VS protocol can organize the most actives at the top of a rank-ordered list (the “early recognition problem”). A VS algorithm that can identify all active molecules but fails to rank them at the top is useless. This is because we only carry out biological evaluation for about 0.1–10% of the molecules in the rank-ordered list⁽⁵²⁾. Some metrics such as the (RIE) and (BEDROC) tackle the issue of “early recognition” in VS⁽⁵⁰⁾. A brief explanation of each of the enrichment descriptors is given below:

4.2.1. Enrichment factor (EF)

The enrichment factor (EF) indicates the capability of the docking software to point out true positives (active ligands) from the total database in comparison with random selection. In general, the following equation is used to estimate EF*

$$EF(\text{subset}) = \frac{\left[\frac{\text{Ligands(selected)}}{N(\text{subset})} \right]}{\left[\frac{\text{Ligands(total)}}{N(\text{total})} \right]} \quad (3)$$

* Where ligands represent the active compounds in the database and N (subset) is the number of the compounds (eg: 100)

obtained by VS software from the total number of compounds present in the database N (total)⁽⁵⁴⁾.

A more simplified form for EF calculation is given by the equation⁽⁵⁵⁾:

$$EF(X\%) = (100/X) * (\text{Fraction of Actives Found}) \quad (4)$$

Where X is a given fraction of the ranked data set for which enrichment is calculated and “fraction of actives” is the ratio between the retrieved active compounds by VS software to the total number of active compounds that exist in the original database (Ligands selected)/ (Ligands total).

The enrichment factor is evaluated at different fractions of the ranked database. Early enrichment is usually calculated for the top 1% of the ranked hit list. Besides, late-stage database screening is addressed through the calculation of the enrichment factor at 20% of the ranked database (EF20). The higher the value of EF is better, but it should be noted that the maximum attainable value is obtained when the VS software retrieves all the active compounds present in the database; (Fraction of Actives Found=1). The maximum values for EF(2%), EF(5%), and EF(10%) are 50, 20, and 10, respectively⁽⁵⁶⁾. The EF is highly dependent on the proportion of the actives in the tested database and thus it can only be used to compare various VS workflows when the same database of actives and decoys is used for assessment. EF doesn't represent the ranking ability of VS software since all actives participate equally to the value regardless of their position in the rank-ordered list⁽⁵²⁾.

4.2.2. Receiver Operating Characteristic (ROC)

The ROC approach, originating from signal detection analysis, is widely used in many fields. It involves plotting true-positive fractions (sensitivity) against false-positive ones (1-specificity) for all compounds in a ranked dataset. Plotting the ROC curve

represents both the sensitivity and specificity of the VS protocol^(54, 57, 58).

ROC curves close to the upper-left corner indicate better performance of VS workflow in discriminating actives from decoys. Theoretically, in ultimate distributions where there is a clear distinction between active compounds and inactive ones, the curve is aligned vertically to the upper-left corner. A diagonal ascending from the origin to the upper right would indicate a random classification of the compounds signifying poor VS protocol (**Fig. 3**)^(52, 59).

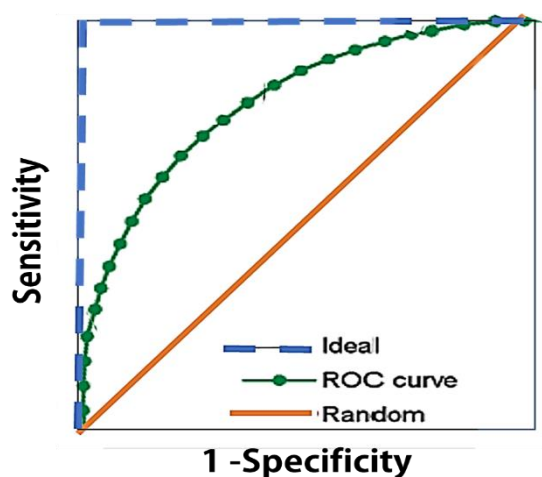


Fig. 3. ROC curves (Sensitivity as a function of 1 -Specificity) in 3 different cases representing ideal, practical (ROC curve), and random distributions of actives and decoys

4.2.3. Area Under the Curve (AUC)

It is a practical way of interpreting the results of ROC curves. The value of the (AUC) is obtained by summing the areas of all the rectangles formed by the Se and 1-Sp values for the different thresholds. AUC is bounded between zero and 1, where value of zero indicates poor performance while 1 indicates efficient VS protocols. Generally speaking, higher values of AUC denote the success of VS process in identifying actives from inactive molecules in the database. The ideal

active and decoy distributions should have a value of 1, while a random distribution results in an AUC value of 0.5.

On the other hand, AUC value less than 0.5 indicates a VS method's unfavourable scenario, where higher scores are given to decoys rather than to actives. ROC and AUC metrics do not rely on the ratio of active compounds to decoys in a dataset^(52, 57).

4.2.4. Robust Initial Enhancement (RIE)

The RIE descriptor describes how many times the distribution of the ranks for active molecules caused by a VS software is superior to a random rank distribution⁽⁵²⁾. It uses a continuously decreasing exponential weight α as a function of rank, where the active molecule located at the beginning of the ordered list gets a weight of approximately 1, and decreasing weights are given for increasing ranks of the actives. This function for the retrieved actives is as follows: $S = \sum_{i=1}^{\text{actives}} \exp\left(-\frac{\text{rank}(i)}{\alpha}\right)$ (5)

The value [S] denotes the total of all weights for all active molecules. This value is divided by the mean sum (S), which is derived from calculations where the active molecules get randomly selected ranks to obtain the final RIE descriptor according to the equation⁽⁵⁰⁾

$$\text{RIE} = \frac{S}{(S)} \quad (6)$$

RIE is similar in meaning to EF and hence, $1/\alpha$ for RIE resembles the X for EF indicating the fraction of the hit list investigated (eg: $\alpha=20$ and 100 meaning the first 20 and 100 retrieved compounds, respectively).

A VS technique that is able to score more active molecules at the top of a rank-ordered list than a random distribution is generally indicated by RIE values larger than 1. Values of 1 indicate a random rank distribution of the active molecules. Therefore, RIE is sensitive to early recognition problem. RIE metric is highly affected by the ratio of actives to in actives and an adjustable parameter, the exponential weight^(52, 57).

4.2.5. Boltzmann- Enhanced Discrimination of ROC (BEDROC)

It is a generalized AUC descriptor that contains a decreasing exponential weighting function that emphasizes active molecules ranked at the top of the ordered list. In this context, this descriptor addresses the early recognition problem in addition, it measures the overall performance of VS process not at a preselected threshold as other descriptors⁽⁵²⁾. This function is derived from RIE metric and hence has an α value (tuning parameter) similar in meaning to RIE. BEDROC value is between [0, 1] it represents the probability that an active is ranked by VS method before a randomly selected compound from a hypothetical exponential probability distribution function with parameter α ⁽⁵⁷⁾. The value of this tuning parameter α concludes what part of the curve is “early”. The suggested α value for BEDROC is 20, which indicates that 80% of the BEDROC value is contributed by the first 8% of the relative rank. A value of $\alpha=160$ is widely used and corresponds to 1% of the actives in the rank-ordered list. It should be noted that BEDROC scores are only comparable when the same α values are used⁽⁵⁰⁾.

This metric can be considered as a “VS usefulness scale” with probability meaning like ROC curve and thus BEDROC values more than 0.5 account for a good VS workflow able to retrieve actives with high ranks within the screened dataset⁽⁵²⁾.

5. Benchmarking dataset; a requirement for enrichment calculations

Retrospective assessment of VS approaches is carried out through examination of test set enrichment or benchmarking datasets. These sets consist of known actives and putative inactives, often referred to as decoys. The decoys selection process is very crucial in evaluating enrichment factors in docking screens⁽⁶⁰⁾. Thus, on designing the decoys, they should resemble the physical properties

of the active ligands while being chemically distinct from the ligands to be non- binders of the target structure thus leading to biased enrichments with artificially good results⁽⁶¹⁻⁶³⁾.

Investigators have put together sets of ligands and suspected decoys for numerous targets leading to the development of numerous benchmarking datasets⁽⁶⁰⁾. The Directory of Useful Decoys (DUD) is a widely used benchmarking dataset that comprises extensive decoy sets for a variety of protein targets and active ligands while maintaining physical similarity. This dataset is widely used for the evaluation of docking methods. An improvement of this dataset has been made, leading to a database named, enhanced DUD (DUD-E). the (DUD-E) is based on the intersection of the databases: ChEMBL for ligand affinities, ZINC for inactive molecules, and the RCSB-PDB database for target structures selection. This dataset provides active ligands and decoy sets for totalling 102 proteins. One can get the entire DUD-E benchmarking set for free at <http://dude.docking.org>⁽⁶⁴⁾. Moreover, with just a list of ligands, the online protocol used to create the decoys for DUD-E can be used to create decoys for any target, making it possible to explore new targets.

6. Pharmacophore modelling (Ligand-based virtual screening)

The concept of a pharmacophore has been applied in many software packages and has many applications in the drug discovery and development process⁽⁶⁵⁾. Building a pharmacophore model can be made through one of two approaches: Ligand-based pharmacophore modelling and Structure-based pharmacophore modelling (E-pharmacophore). Ligand-based modelling relies on knowing the 3D structures of active ligands against a specific target that are superimposed to extract common chemical features responsible for bioactivity. This approach does not require knowledge of the

target protein structure and hence, could be used for pharmacophore model generation for target structures not fully resolved by NMR or X-ray crystallographic techniques⁽¹⁵⁾. Structure-based pharmacophore modeling relies mainly on the 3D structure of a macromolecular target either free or with its bound ligand. This pharmacophoric model is derived based on analyzing complementary chemical features of the target active site and their spatial relationships^(15, 65). Pharmacophore models are designed to provide relevant information for drug design by easily interpreting the locations of functional groups involved in ligand-target interactions and the various forms of non-covalent bonding⁽⁶⁶⁾.

Although their representation varies among different software, these models are commonly represented by spatial arrangement (Fig. 4) that contain the pharmacophoric points (fragments, chemical features) and the geometrical constraints connecting them (distances, torsions and 3D-coordinate location constraints as excluded volumes)⁽⁶⁷⁾.

Regions of “forbidden” space that the active molecule should not occupy to avoid a steric clash with the target are represented in pharmacophore models as “variably sized spheres excluded volumes”. The hits obtained following pharmacophore model filtration must fit the interaction features defined by the model and fit within the region.

7. 3D-Quantitative structure-activity relationship (3D-QSAR) models

QSAR has been known for decades, simply, it's a process for developing mathematical or computer models to measure the relationships between chemical substances' physicochemical characteristics and biological activities through the use of chemometric techniques. These QSAR models should be reliable and statistically significant for the prediction of the activities

of new chemical entities⁽⁶⁹⁾. To build 3D-QSAR model, first, experimental activity data for enzyme inhibitors or receptor ligands should be available. This is followed by superimposition or alignment of the 3D structures of these molecules to identify and determine molecular descriptors of these molecules. Then, finding the correlations between these descriptors and the biological activity and finally, testing the statistical stability and predictive power of the developed 3D-QSAR model⁽⁷⁰⁾.

Different approaches are used for 3D-QSAR model development according to the software applied. In most cases, developing a pharmacophore model comes as a preceding step to the generation of 3D-QSAR model as a superimposition process^(69, 71).

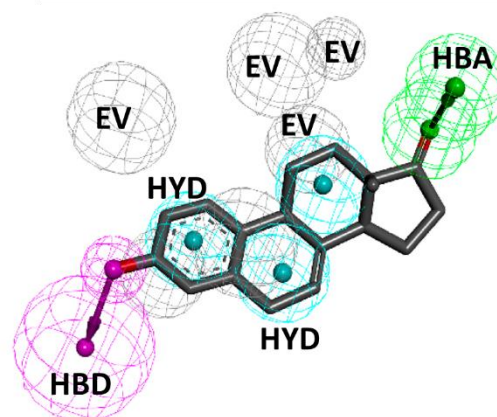


Fig. 4: Representation of pharmacophore model with the active compound aligned where HBA; hydrogen bond acceptor, HBD; hydrogen bond donor, HYD; hydrophobic, EV; excluded volume. Figure adapted from Kaserer et al.⁽⁶⁸⁾, which was published under the Creative Commons by Attribution (CC-BY) license

8. Training set specification for building pharmacophore and 3D-QSAR models

The development of either pharmacophore / 3D-QSAR models for a specific target using (The ligand-based approach) is often

accomplished by taking a series of molecules' 3D structures and extracting common chemical characteristics from them interacting with the macromolecular target (training set) ⁽⁷²⁾. Selecting training set compounds is a very challenging task. The size, chemical variety, and kind of ligand molecules all have a significant influence on the pharmacophore model produced. It has been demonstrated that different pharmacophore models could be generated for the same molecular target using the same program when altered training set compounds are used ⁽⁷³⁾. For pharmacophore models, selecting the training set is a simple and flexible process. Compounds in the training set could be as minimum as two compounds with no defined maximum number composing the training set ⁽⁷⁴⁾.

In the case of 3D-QSAR, the training set implemented in pharmacophore model construction could serve as a set for 3D-QSAR model generation provided that the compounds possess the following criteria ⁽⁶⁹⁾;

- Compounds with identical binding modes share the same mechanism of action.
- The biological activity of the chosen compounds should be assayed using the same experimental study and the results expressed using the same unit of measurements (preferably data obtained from a single source).
- Biological activity data should cover a wide range (4 orders of magnitude is recommended), where several ligands should be included from all categories from the most active to the least ones even the inactives are included ⁽⁷⁵⁾.

The minimum number of training set compounds should be 16 to ensure a statistically significant model, while no limits exist for the maximum number used ^(74, 76).

To obtain compounds for training sets, many public databases are available such as ChEMBL and Drugbank that include activity data of compounds. In addition, ToxCast,

Tox21, and PubChem Bioassay provide data for both active and inactive molecules. These databases could be used for extracting such compounds with their activity data besides the literature review ⁽⁷⁴⁾.

9. 3D-QSAR model validation

QSAR models must be validated to confirm their reliability; and the capability of the model to predict compounds' activity outside the training set. There are two types of validation used for 3D-QSAR models: (a) internal validation, and (b) external validation ^(69, 77).

9.1. Internal validation

It is performed on the training set used for developing the 3D-QSAR model. The most common method used for this purpose is "leave-one-out" (LOO) - cross-validation approach ⁽⁷⁰⁾. This is done via testing the ability of the model composed from the training dataset to predict one of the observations previously eliminated from it. This process is repeated so that all of the observations play the role of a test compound ^(78, 79).

The "leave-n-out" cross-validation method is considered more reliable and robust where a model constructed from the remaining compounds is used to predict the activity of small subsets of the training set that are temporarily held out ⁽⁷¹⁾. The predictive power of the model is then assessed by calculating the R^2 of the training set. R^2 (R-squared; coefficient of determination) shows how resilient the model is to modifications in the training set. It has a maximum value and is computed between the leave-n-out predictions and the predictions made by the model constructed using the entire training set of 1 ⁽⁷⁹⁾.

9.2. External validation

External validation involves predicting the activity of new compounds not included in the original true test set. This method is considered far superior to the internal cross-validation process ^(78, 79). The statistical

parameters used to judge external validation are Q^2 (predictive squared correlation coefficient) and, RMSE (root mean square error) ^(77, 80).

It should be noted that many different statistical parameters are used for 3D-QSAR model validation. Some of the commonly used parameters for the training set are (SD) Standard deviation of the regression, R^2 for the regression, F (Fischer statistic or variance ratio) which can indicate the level of statistical significance of the regression model where higher values indicate more significant regression, and its level of significance (P-value). For test set: (RMSE) Root-mean-square error, (Q^2) for the predicted activities, and (Pearson-R) for the correlation between the predicted and observed activity for the test set ⁽⁶⁹⁾. Moreover, plotting the experimental activities against the predicted ones for the training and test sets and calculating their R^2 is a further assessment of the reliability of the developed 3D-QSAR models. R^2 values near 1 indicate models with more predictive ability of the model especially for molecules outside the training set ⁽³⁷⁾.

10. Advantages and limitations of virtual screening techniques

Virtual screening techniques were proved a powerful strategy for speeding up the drug discovery process. Screening gigantic databases in a timewise and cost-effective manner is a breakthrough of VS techniques. Compared to high throughput screening (HTS), VS is performed in-silico which doesn't necessitate the existence of all the compounds. This led to the identification of active hits with reduced costs ^(40, 81). Nowadays, a large number of docking softwares and pharmacophore modelling programs exist. This facilitates the identification of active drug candidates with nM potency. Moreover, VS can be easily utilized as a screening tool in finding new scaffolds with promising activities. A

literature survey showed many success stories for the discovery of potent drugs through VS approaches. Captopril, indinavir, saquinavir, and others are a few examples ⁽¹⁴⁾. Unfortunately, VS techniques suffer from some limitations. This is mainly due to the fact that living biological systems are very complex and their simulation using computational methods is impractical. Furthermore, the scoring functions used in VS usually lead to a high rate of false positives and/or false negatives results. VS protocols are highly dependent on the database used; thus the results can vary greatly even for the same target.

Docking simulation algorithms suffer from the limited flexibility of the protein or enzyme structure while offering high flexibility of the screened ligands. In living systems, both receptor and target molecules are highly flexible in solution, thus docking using limited protein flexibilities may lead to the wrong results. Some algorithms don't incorporate the effect of water molecules or other solvents that have a fundamental role in cellular systems.

The generation of pharmacophore models doesn't take into account the dynamic nature of the receptor and ligand due to their flexibility in living cells. This leads to models lacking important features that could participate in predicting potential hits ⁽⁸²⁾. A major limitation of the pharmacophore model is that they include the chemical nature of the ligands and are missing variables that can influence the binding of active compounds to the specific target. Features such as chemical solubility, cell metabolism, membrane permeability, and others that greatly affect ligand binding to the protein active site are not addressed during pharmacophore model generation ⁽⁸³⁾.

11. Cloud computing and deep learning; the future of VS techniques

Over the past few years and with the achievements in computer technologies, the

virtual screening process has been enhanced especially on the scalability level⁽⁹⁾. Virtual screening could work efficiently with a few hundred thousand and a few million molecules. Recently, with the availability of databases containing a billion molecules or more, virtual screening could be performed for such gigantic databases utilizing cloud computing infrastructure⁽⁸⁴⁾. Cloud computing infrastructure provides virtual machines and millions of processors and CPUs for a relatively low cost at any time. It is a means to integrate different workflows and tools provided through service providers. Clouding computing can be performed using many service providers such as Amazon Web Services and Google Cloud. This allows the conduction of virtual screens for a massive number of molecules at an affordable cost. Cloud computing can replace local compute cluster used in research facilities saving the expenses of hardware, maintenance, and others⁽⁸⁵⁻⁸⁷⁾. Another major advantage of VS techniques is the way the data is being processed to provide meaningful results. A major achievement is the use of deep learning methods complementing the traditionally used machine learning methods. Deep learning is classified as a subset of machine learning (ML) where the latter aims at learning and implementing tasks that can be used to predict certain parameter (s). ML is utilized in VS through training models using datasets of both active and inactive compounds then, algorithms are generated based on pairs of inputs and outputs, This can then be utilized to forecast a compound's activity⁽⁸⁸⁾. In deep learning (DL), a network is created by processing and transforming the input data into several hidden layers. These hidden layers consist of a linear component and a non-linear component, which is referred to as the activation function.^(89, 90). The information-driven through these networks results in obtaining highly flexible predictive models. Deep neural networks

(DNNs) are one example of the deep learning methods applied in VS⁽⁹¹⁾. These networks transform the input data into more complex features and thus could have great performance in describing the complex interactions existing between molecules and biological targets⁽⁸⁸⁾. Their application in virtual screening resulted in many advantages, for instance, identifying relationships between multiple targets. Moreover, deep learning methods allow the prediction of hit compounds for targets having less training samples based on the shared hidden units among the targets⁽⁹²⁾. Deep learning methods adopted in VS have a great impact on improving the overall performance of VS. They show very promising results in comparison with the machine learning and docking classical approaches.

12. Conclusions

Both molecular docking and pharmacophore modelling stand out as an effective computational technique highly applicable during the drug discovery process. Being predictive techniques performed on virtual levels, their validation is of great importance to obtain reliable results. Understanding the basic concepts of these validations allows researchers to have more accurate insights into the results of VS workflow. It also helps explain the reliability of the obtained results. Validation of molecular docking is performed before conducting the screening. meanwhile, that of 3D-QSAR is done after model development to judge the stability and reliability of the model. Despite the wide applications of VS strategies in the drug discovery process, they suffer from many limitations. In particular, a shortage of these techniques to simulate the environment in living cells. It is strongly recommended to combine different VS strategies (SBVS and LBVS) to complement each other.

The field of VS has rapidly evolved with the recent development in computing power and

machine learning techniques. Cloud computing and deep learning are among the newly evolved technologies contributing to the development of VS techniques. This requires the merge between scientific knowledge on molecular interactions levels with the enormous technological advances to provide meaningful results applicable during the drug discovery process.

Authors Contributions

H F., R I. introduced the concept. H F., R I., and R M. wrote the drafts, participated in reviewing them, and made the final approval for publication. All authors contributed equally to data gathering.

Conflict of interest

The authors declare no conflict of interest.

13. References

- (1) Petersen F, Amstutz R. *Natural Compounds as Drugs*; Birkhäuser Basel; 2007.
- (2) Potterat O, Hamburger M. Natural products in drug discovery-concepts and approaches for tracking bioactivity. *Current Organic Chemistry*. 2006;10(8):899-920.
- (3) Lin X, Li X, Lin X. A Review on Applications of Computational Methods in Drug Screening and Design. *Molecules*. 2020;25(6):1375.
- (4) Shaker B, Ahmad S, Lee J, Jung C, Na D. In silico methods and tools for drug discovery. *Computers in Biology and Medicine*. 2021;137:104851.
- (5) Brogi S, Ramalho TC, Kuca K, Medina-Franco JL, Valko M. In silico methods for drug design and discovery. *Frontiers Media SA*; 2020. p. 612.
- (6) Ferreira LLG, Andricopulo AD. ADMET modeling approaches in drug discovery. *Drug Discovery Today*. 2019;24(5):1157-65.
- (7) Ma D-L, Chan DS-H, Leung C-H. Molecular docking for virtual screening of natural product databases. *Chemical science*. 2011;2(9):1656-65.
- (8) Shen J, Xu X, Cheng F, Liu H, Luo X, Chen K, et al. Virtual screening on natural products for discovering active compounds and target information. *Curr Med Chem*. 2003;10(21):2327-42.
- (9) Walters WP, Wang R. New Trends in Virtual Screening. *Journal of Chemical Information and Modeling*. 2020;60(9):4109-11.
- (10) Dhasmana A, Raza S, Jahan R, Lohani M, Arif JM. Chapter 19 - High-Throughput Virtual Screening (HTVS) of Natural Compounds and Exploration of Their Biomolecular Mechanisms: An In Silico Approach. In: Ahmad Khan MS, Ahmad I, Chattopadhyay D, editors. *New Look to Phytomedicine*; Academic Press; 2019. p. 523-48.
- (11) Subhaswaraj P, Siddhardha B. Chapter 11 - Molecular docking and molecular dynamic simulation approaches for drug development and repurposing of drugs for severe acute respiratory syndrome-Coronavirus-2. In: Parihar A, Khan R, Kumar A, Kaushik AK, Gohel H, editors. *Computational Approaches for Novel Therapeutic and Diagnostic Designing to Mitigate SARS-CoV-2 Infection*; Academic Press; 2022. p. 207-46.
- (12) Lavecchia A, Giovanni DC. Virtual Screening Strategies in Drug Discovery: A Critical Review. *Current Medicinal Chemistry*. 2013;20(23):2839-60.
- (13) Jiang L, Rizzo RC. Pharmacophore-Based Similarity Scoring for DOCK. *The Journal of Physical Chemistry B*. 2015;119(3):1083-102.
- (14) Baig MH, Ahmad K, Roy S, Ashraf JM, Adil M, Siddiqui MH, et al. Computer Aided Drug Design: Success and Limitations. *Curr Pharm Des*. 2016;22(5):572-81.
- (15) Yang S-Y. Pharmacophore modeling and applications in drug discovery: challenges and recent advances. *Drug Discovery Today*. 2010;15(11):444-50.
- (16) Brooijmans N, Kuntz ID. Molecular recognition and docking algorithms. *Annual review of biophysics and biomolecular structure*. 2003;32:335-73.
- (17) Sousa SF, Fernandes PA, Ramos MJ. Protein-ligand docking: current status and future challenges. *Proteins*. 2006;65(1):15-26.
- (18) Morris GM, Goodsell DS, Halliday RS, Huey R, Hart WE, Belew RK, et al. Automated docking using a Lamarckian genetic algorithm and an empirical binding free energy function. *Journal of Computational Chemistry*. 1998;19(14):1639-62.
- (19) Ewing TJA, Makino S, Skillman AG, Kuntz ID. DOCK 4.0: Search strategies for automated molecular docking of flexible molecule databases. *Journal of Computer-Aided Molecular Design*. 2001;15(5):411-28.
- (20) Warren GL, Andrews CW, Capelli A-M, Clarke B, LaLonde J, Lambert MH, et al. A Critical Assessment of Docking Programs and Scoring Functions. *Journal of Medicinal Chemistry*. 2006;49(20):5912-31.
- (21) McGann M. FRED and HYBRID docking performance on standardized datasets. *J Comput Aided Mol Des*. 2012;26(8):897-906.

- (22) Friesner RA, Banks JL, Murphy RB, Halgren TA, Klicic JJ, Mainz DT, et al. Glide: a new approach for rapid, accurate docking and scoring. 1. Method and assessment of docking accuracy. *J Med Chem*. 2004;47(7):1739-49.
- (23) Verdonk ML, Cole JC, Hartshorn MJ, Murray CW, Taylor RD. Improved protein-ligand docking using GOLD. *Proteins*. 2003;52(4):609-23.
- (24) Abagyan R, Totrov M, Kuznetsov D. ICM—A new method for protein modeling and design: Applications to docking and structure prediction from the distorted native conformation. *Journal of Computational Chemistry*. 1994;15(5):488-506.
- (25) Schrödinger L. Schrödinger, LLC; New York, NY: 2017. Schrödinger Suite. 2017;2:2017-1.
- (26) Zavodszky MI, Rohatgi A, Van Voorst JR, Yan H, Kuhn LA. Scoring ligand similarity in structure-based virtual screening. *J Mol Recognit*. 2009;22(4):280-92.
- (27) Jain AN. Surflex-Dock 2.1: robust performance from ligand energetic modeling, ring flexibility, and knowledge-based search. *J Comput Aided Mol Des*. 2007;21(5):281-306.
- (28) Thomsen R, Christensen MH. MolDock: A New Technique for High-Accuracy Molecular Docking. *Journal of Medicinal Chemistry*. 2006;49(11):3315-21.
- (29) Chen R, Li L, Weng Z. ZDOCK: an initial-stage protein-docking algorithm. *Proteins*. 2003;52(1):80-7.
- (30) Barnum D, Greene J, Smellie A, Sprague P. Identification of Common Functional Configurations Among Molecules. *Journal of Chemical Information and Computer Sciences*. 1996;36(3):563-71.
- (31) Klebe G, Abraham U, Mietzner T. Molecular Similarity Indices in a Comparative Analysis (CoMSIA) of Drug Molecules to Correlate and Predict Their Biological Activity. *Journal of Medicinal Chemistry*. 1994;37(24):4130-46.
- (32) Martin YC, Bures MG, Danaher EA, DeLazzer J, Lico I, Pavlik PA. A fast new approach to pharmacophore mapping and its application to dopaminergic and benzodiazepine agonists. *Journal of Computer-Aided Molecular Design*. 1993;7(1):83-102.
- (33) Richmond NJ, Abrams CA, Wolohan PRN, Abrahamian E, Willett P, Clark RD. GALAHAD: 1. Pharmacophore identification by hypermolecular alignment of ligands in 3D. *Journal of Computer-Aided Molecular Design*. 2006;20(9):567-87.
- (34) Barnum D, Greene J, Smellie A, Sprague P. Identification of common functional configurations among molecules. *J Chem Inf Comput Sci*. 1996;36(3):563-71.
- (35) Li H, Sutter J, Remy H. Pharmacophore Perception, Development, and Use in Drug Design, ch. HypoGen: An Automated System for Generating 3D Predictive Pharmacophore Models. International University Line. 2000:49-68.
- (36) Wolber G, Langer T. LigandScout: 3-D Pharmacophores Derived from Protein-Bound Ligands and Their Use as Virtual Screening Filters. *Journal of Chemical Information and Modeling*. 2005;45(1):160-9.
- (37) Dixon SL, Smondryev AM, Knoll EH, Rao SN, Shaw DE, Friesner RA. PHASE: a new engine for pharmacophore perception, 3D QSAR model development, and 3D database screening: 1. Methodology and preliminary results. *Journal of Computer-Aided Molecular Design*. 2006;20(10):647-71.
- (38) Joseph-McCarthy D, Alvarez JC. Automated generation of MCSS-derived pharmacophoric DOCK site points for searching multiconformation databases. *Proteins: Structure, Function, and Bioinformatics*. 2003;51(2):189-202.
- (39) Prachayasittikul V, Worachartcheewan A, Shoombuatong W, Songtawee N, Simeon S, Prachayasittikul V, et al. Computer-Aided Drug Design of Bioactive Natural Products. *Curr Top Med Chem*. 2015;15(18):1780-800.
- (40) Lionta E, Spyrou G, Vassilatis DK, Cournia Z. Structure-based virtual screening for drug discovery: principles, applications and recent advances. *Curr Top Med Chem*. 2014;14(16):1923-38.
- (41) Dias R, de Azevedo WF. Molecular docking algorithms. *Current drug targets*. 2008;9(12):1040-7.
- (42) Meng XY, Zhang HX, Mezei M, Cui M. Molecular docking: a powerful approach for structure-based drug discovery. *Curr Comput Aided Drug Des*. 2011;7(2):146-57.
- (43) Brünger AT, Nilges M. Computational challenges for macromolecular structure determination by X-ray crystallography and solution NMR spectroscopy. *Quarterly Reviews of Biophysics*. 2009;26(1):49-125.
- (44) Cavasotto CN, Phatak SS. Homology modeling in drug discovery: current trends and applications. *Drug Discovery Today*. 2009;14(13):676-83.
- (45) Jacob RB, Andersen T, McDougal OM. Accessible high-throughput virtual screening molecular docking software for students and

- educators. *PLoS computational biology*. 2012;8(5):e1002499.
- (46) Macarron R, Banks MN, Bojanic D, Burns DJ, Cirovic DA, Garyantes T, et al. Impact of high-throughput screening in biomedical research. *Nature Reviews Drug Discovery*. 2011;10(3):188-95.
- (47) Mukherjee S, Balius TE, Rizzo RC. Docking Validation Resources: Protein Family and Ligand Flexibility Experiments. *Journal of Chemical Information and Modeling*. 2010;50(11):1986-2000.
- (48) Hevener KE, Zhao W, Ball DM, Babaoglu K, Qi J, White SW, et al. Validation of Molecular Docking Programs for Virtual Screening against Dihydropteroate Synthase. *Journal of Chemical Information and Modeling*. 2009;49(2):444-60.
- (49) Rueda M, Abagyan R. Best Practices in Docking and Activity Prediction. *bioRxiv*. 2016:039446.
- (50) Truchon J-F, Bayly CI. Evaluating Virtual Screening Methods: Good and Bad Metrics for the “Early Recognition” Problem. *Journal of Chemical Information and Modeling*. 2007;47(2):488-508.
- (51) Yuriev E, Agostino M, Ramsland PA. Challenges and advances in computational docking: 2009 in review. *J Mol Recognit*. 2011;24(2):149-64.
- (52) Kirchmair J, Markt P, Distinto S, Wolber G, Langer T. Evaluation of the performance of 3D virtual screening protocols: RMSD comparisons, enrichment assessments, and decoy selection--what can we learn from earlier mistakes? *J Comput Aided Mol Des*. 2008;22(3-4):213-28.
- (53) Zhao W, Hevener KE, White SW, Lee RE, Boyett JM. A statistical framework to evaluate virtual screening. *BMC bioinformatics [Internet]*. 2009 2009/07//; 10:[225 p.].
- (54) Huang N, Shoichet BK, Irwin JJ. Benchmarking sets for molecular docking. *J Med Chem*. 2006;49(23):6789-801.
- (55) Nicholls A. What do we know and when do we know it? *Journal of Computer-Aided Molecular Design*. 2008;22(3):239-55.
- (56) Friesner RA, Banks JL, Murphy RB, Halgren TA, Klicic JJ, Mainz DT, et al. Glide: A New Approach for Rapid, Accurate Docking and Scoring. 1. Method and Assessment of Docking Accuracy. *Journal of Medicinal Chemistry*. 2004;47(7):1739-49.
- (57) Riniker S, Landrum GA. Open-source platform to benchmark fingerprints for ligand-based virtual screening. *Journal of Cheminformatics*. 2013;5(1):26.
- (58) Empereur-Mot C, Guillemain H, Latouche A, Zagury J-F, Viallon V, Montes M. Predictiveness curves in virtual screening. *Journal of cheminformatics [Internet]*. 2015 2015; 7:[52 p.].
- (59) Triballeau N, Acher F, Brabet I, Pin J-P, Bertrand H-O. Virtual Screening Workflow Development Guided by the “Receiver Operating Characteristic” Curve Approach. Application to High-Throughput Docking on Metabotropic Glutamate Receptor Subtype 4. *Journal of Medicinal Chemistry*. 2005;48(7):2534-47.
- (60) Lagarde N, Zagury J-F, Montes M. Benchmarking Data Sets for the Evaluation of Virtual Ligand Screening Methods: Review and Perspectives. *Journal of Chemical Information and Modeling*. 2015;55(7):1297-307.
- (61) Wallach I, Lilien R. Virtual Decoy Sets for Molecular Docking Benchmarks. *Journal of Chemical Information and Modeling*. 2011;51(2):196-202.
- (62) Vogel SM, Bauer MR, Boeckler FM. DEKOIS: Demanding Evaluation Kits for Objective in Silico Screening — A Versatile Tool for Benchmarking Docking Programs and Scoring Functions. *Journal of Chemical Information and Modeling*. 2011;51(10):2650-65.
- (63) Rohrer SG, Baumann K. Maximum Unbiased Validation (MUV) Data Sets for Virtual Screening Based on PubChem Bioactivity Data. *Journal of Chemical Information and Modeling*. 2009;49(2):169-84.
- (64) Mysinger MM, Carchia M, Irwin JJ, Shoichet BK. Directory of Useful Decoys, Enhanced (DUD-E): Better Ligands and Decoys for Better Benchmarking. *Journal of Medicinal Chemistry*. 2012;55(14):6582-94.
- (65) Caporuscio F, Tafi A. Pharmacophore modelling: a forty year old approach and its modern synergies. *Current medicinal chemistry*. 2011;18(17):2543-53.
- (66) Seidel T, Bryant SD, Ibis G, Poli G, Langer T. 3D pharmacophore modeling techniques in computer-aided molecular design using LigandScout. *Tutorials Chemoinform*. 2017;281:279-309.
- (67) Poptodorov K, Luu T, Hoffmann RD. Pharmacophore model generation software tools. 2006.
- (68) Kaserer T, Beck KR, Akram M, Odermatt A, Schuster D. Pharmacophore Models and Pharmacophore-Based Virtual Screening: Concepts and Applications Exemplified on Hydroxysteroid Dehydrogenases. *Molecules*. 2015;20(12):22799-832.
- (69) Verma J, Khedkar VM, Coutinho EC. 3D-QSAR in drug design--a review. *Current topics in medicinal chemistry*. 2010;10(1):95-115.

- (70) Acharya C, Coop A, Polli JE, Mackerell AD. Recent advances in ligand-based drug design: relevance and utility of the conformationally sampled pharmacophore approach. *Current computer-aided drug design*. 2011;7(1):10-22.
- (71) Sippl W. Application of Structure-Based Alignment Methods for 3D QSAR Analyses. *Pharmacophores and Pharmacophore Searches*. 2006;32:223-49.
- (72) Mannhold R, Kubinyi H, Folkers G. *Pharmacophores and pharmacophore searches*: John Wiley & Sons; 2006.
- (73) Wolber G, Seidel T, Bendix F, Langer T. Molecule-pharmacophore superpositioning and pattern matching in computational drug design. *Drug Discovery Today*. 2008;13(1):23-9.
- (74) Kaserer T, Beck KR, Akram M, Odermatt A, Schuster D. *Pharmacophore Models and Pharmacophore-Based Virtual Screening: Concepts and Applications Exemplified on Hydroxysteroid Dehydrogenases*. *Molecules* (Basel, Switzerland) [Internet]. 2015 2015/12//; 20(12):[22799-832 pp.]
- (75) Krovat EM, Langer T. Non-Peptide Angiotensin II Receptor Antagonists: Chemical Feature Based Pharmacophore Identification. *Journal of Medicinal Chemistry*. 2003;46(5):716-26.
- (76) Evans DA, Doman TN, Thorner DA, Bodkin MJ. 3D QSAR Methods: Phase and Catalyst Compared. *Journal of Chemical Information and Modeling*. 2007;47(3):1248-57.
- (77) Kiralj R, Ferreira MMC. Basic validation procedures for regression models in QSAR and QSPR studies: theory and application. *Journal of the Brazilian Chemical Society*. 2009;20:770-87.
- (78) Roy K, Ambure P, Aher RB. How important is to detect systematic error in predictions and understand statistical applicability domain of QSAR models? *Chemometrics and Intelligent Laboratory Systems*. 2017;162:44-54.
- (79) Roy K, Mitra I. On various metrics used for validation of predictive QSAR models with applications in virtual screening and focused library design. *Comb Chem High Throughput Screen*. 2011;14(6):450-74.
- (80) Consonni V, Ballabio D, Todeschini R. Comments on the Definition of the Q2 Parameter for QSAR Validation. *Journal of Chemical Information and Modeling*. 2009;49(7):1669-78.
- (81) Reddy AS, Pati SP, Kumar PP, Pradeep HN, Sastry GN. Virtual Screening in Drug Discovery - A Computational Perspective. *Current Protein and Peptide Science*. 2007;8(4):329-51.
- (82) Qing X, Yin Lee X, De Raeymaeker J, Rh Tame J, Yj Zhang K, De Maeyer M, et al. Pharmacophore modeling: advances, limitations, and current utility in drug discovery. *Journal of Receptor, Ligand and Channel Research*. 2014;7(null):81-92.
- (83) Giordano D, Biancanello C, Argenio MA, Facchiano A. *Drug Design by Pharmacophore and Virtual Screening Approach*. *Pharmaceuticals* (Basel). 2022;15(5).
- (84) Haga JH, Ichikawa K, Date S. *Virtual Screening Techniques and Current Computational Infrastructures*. *Curr Pharm Des*. 2016;22(23):3576-84.
- (85) Olğaç A, Türe A, Olğaç S, Möller S. Cloud-based high throughput virtual screening in novel drug discovery. *High-Performance Modelling and Simulation for Big Data Applications: Selected Results of the COST Action IC1406 cHiPSet*: Springer International Publishing Cham; 2019. p. 250-78.
- (86) Singh N, Sharma U, Bairagee D, Jain N. *Cloud Application in Drug Development. Bioinformatics Tools and Big Data Analytics for Patient Care*: Chapman and Hall/CRC; 2022. p. 181-200.
- (87) Puertas-Martín S, Banegas-Luna AJ, Paredes-Ramos M, Redondo JL, Ortigosa PM, Brovarets OhO, et al. Is high performance computing a requirement for novel drug discovery and how will this impact academic efforts? *Expert Opinion on Drug Discovery*. 2020;15(9):981-5.
- (88) Carpenter KA, Cohen DS, Jarrell JT, Huang X. Deep learning and virtual drug screening. *Future Med Chem*. 2018;10(21):2557-67.
- (89) Kimber TB, Chen Y, Volkamer A. *Deep Learning in Virtual Screening: Recent Applications and Developments*. *Int J Mol Sci*. 2021;22(9).
- (90) Oliveira TAd, Silva MPd, Maia EHB, Silva AMd, Taranto AG. *Virtual Screening Algorithms in Drug Discovery: A Review Focused on Machine and Deep Learning Methods. Drugs and Drug Candidates*. 2023;2(2):311-34.
- (91) Chien J-T. Chapter 7 - Deep Neural Network. In: Chien J-T, editor. *Source Separation and Machine Learning*: Academic Press; 2019. p. 259-320.
- (92) Rifaioglu AS, Atas H, Martin MJ, Cetin-Atalay R, Atalay V, Doğan T. Recent applications of deep learning and machine intelligence on in silico drug discovery: methods, tools and databases. *Brief Bioinform*. 2019;20(5):1878-912.