# GALAS Modeling Methodology Applications in the Prediction of the Drug Safety Related Properties

**Andrius Sazonovas[1], Remigijus Didziapetris[1], Justas Dapkunas[1,2], Liutauras Juska[1,2], Pranas Japertas[1]**

[1] ACD/Labs, Inc., A.Mickeviciaus g. 29, LT-08117 Vilnius, Lithuania,
[2] Department of Biochemistry and Biophysics, Vilnius University, M.K.Ciurlionio g. 21/27, LT-03101 Vilnius, Lithuania.

Advanced Chemistry Development

ACD/Labs

## INTRODUCTION

Early computational evaluation of drug candidate properties related to its pharmaceutical safety (such as hERG inhibition induced cardiotoxicity or CYP3A4 inhibition responsible various unwanted drug-drug interactions) is becoming increasingly important in the drug discovery process. Yet, every model, no matter what data, descriptors or modeling techniques used to build it, has a certain applicability domain, beyond which, the quality of predictions becomes highly questionable. This reality is one of the fundamental issues concerning the effective use of third-party predictive algorithms in industry. The simple reason for this is that literature based training sets rarely cover the specific part of the chemical space that 'in-house' projects are focused on. Discrepancies between 'in-house' experimental protocols and methods used to measure properties for compounds in publicly available sources further affect the quality of resulting *in silico* predictions. Therefore the need has long existed for a method that would allow any company to effectively assess the Applicability Domain of any third-party model and to tailor it to its specific needs using proprietary 'in-house' data.

## GALAS MODEL METHODOLOGY AND RELIABILITY INDEX

Addressing aforementioned issue, a GALAS (Global, Adjusted Locally According to Similarity) model concept has been developed providing a novel solution to this problem. Each GALAS model consists of the following parts:

- Structure based QSAR/QSPR for the prediction of the property of interest – (i.e., baseline model)
- User defined data set with experimental values for the property of interest – (i.e., Self-training Library)
- Special similarity based routine which identifies the most similar compounds contained in the Self-training Library and considering their experimental values calculates systematic deviations produced by the baseline QSAR/QSPR for each submitted molecule – (i.e., training engine)

The result is a prediction that is corrected according to the experimental values for the most similar compounds present in the user defined Self-training Library covering the part of the chemical space not initially included in the training set.
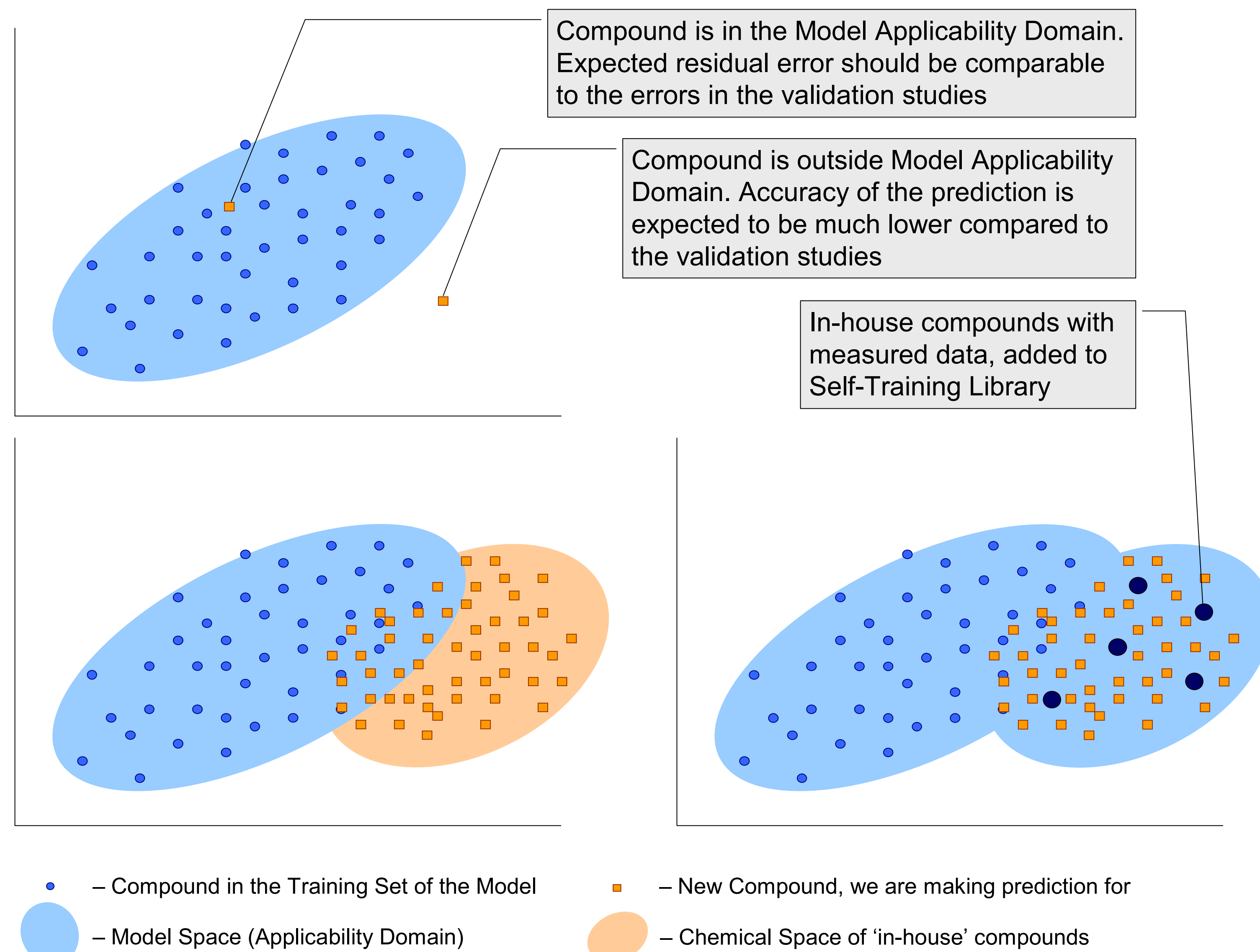


Compound is in the Model Applicability Domain. Expected residual error should be comparable to the errors in the validation studies

Compound is outside Model Applicability Domain. Accuracy of the prediction is expected to be much lower compared to the validation studies

In-house compounds with measured data, added to Self-Training Library

- – Compound in the Training Set of the Model
- – New Compound, we are making prediction for
- – Model Space (Applicability Domain)
- – Chemical Space of 'in-house' compounds

**FIGURE 1.** Illustration of the Model Applicability Domain, and its expansion using GALAS modeling method

In addition, GALAS modeling methodology allows quantitative assessment of the prediction reliability. This information is contained in the developed Reliability Index (RI) that can provide values in the range [0; 1]. Lower values suggest a compound being further from the Model Applicability Domain and the prediction less reliable, on the other hand, high RI values indicate an increasing confidence about the quality of the prediction. Estimation of the Reliability Index takes into account the following aspects:

- Similarity of the tested compound to the training set – no reliable predictions can be made if we have no similar compounds in the training set.
- Consistency of the experimental values for similar compounds – Even when similar compounds are present in the dataset the quality of prediction could be lower if that data is inconsistent.
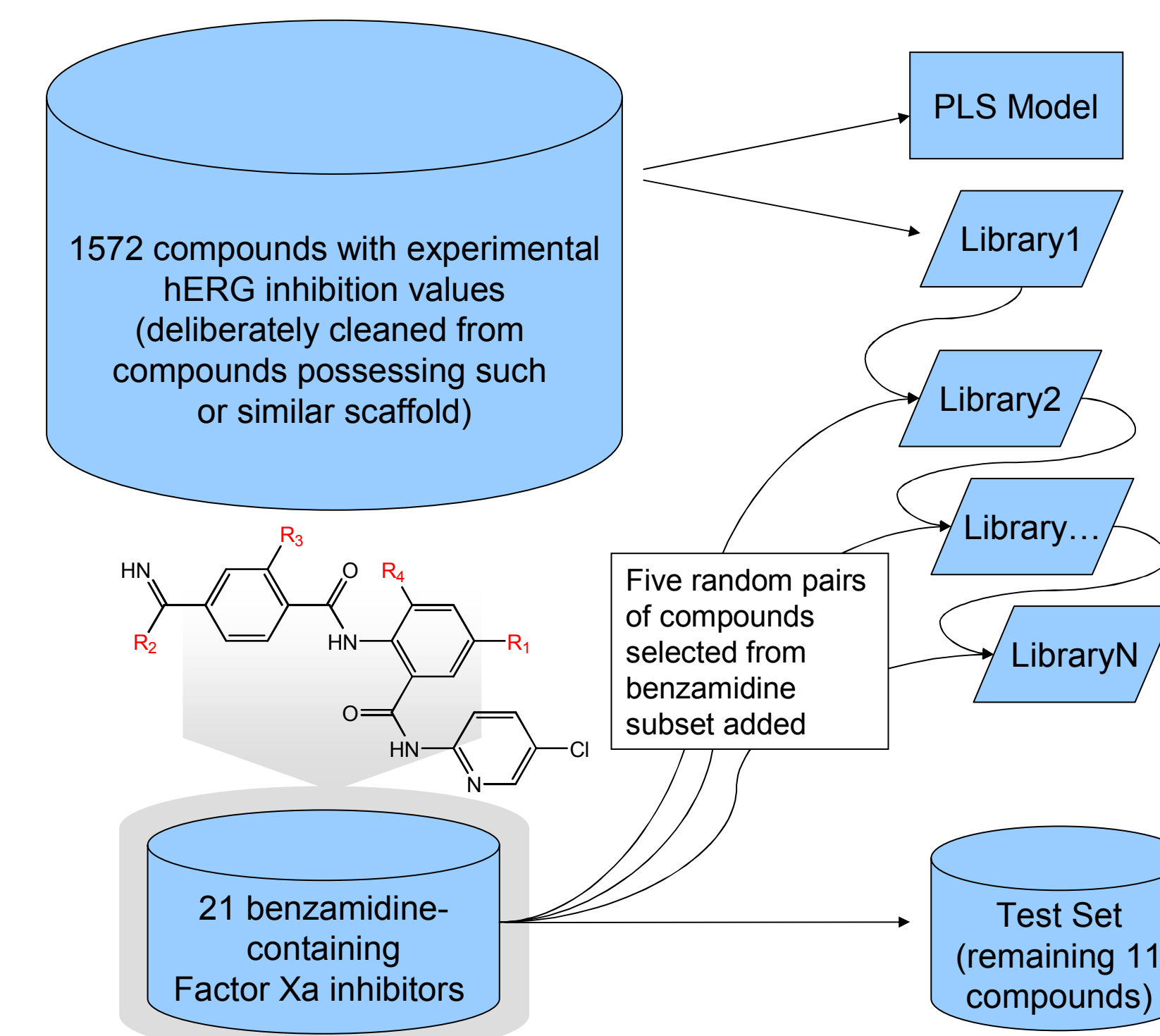
## COPING WITH COMPLETELY NEW CHEMICAL FEATURES – AN EXAMPLE SCENARIO WITH HERG INHIBITION PREDICTION

The objectives of this scenario of the GALAS modeling methodology validation were as follows:

- Demonstrate that a GALAS model can be trained to the completely new chemical features absent in the original training set
- Demonstrate that quite small number of compounds with experimental data is sufficient for such purpose

As outlined in the scheme of the preparation steps, the compound class mimicking a new drug development project in this virtual experiment was benzamidine-containing Factor Xa (thrombokinase) inhibitors [1].

Initially, predictions for all compounds are inconclusive as no similar compounds are present in the library (indicated by the low RI values). However, after just a couple of additions, predictions of sufficient reliability start to appear. When 6 compounds of the same class are added to the library, most calculated values become reliable, and when all 10 compounds are added, 10 of 11 test set molecules are confidently predicted as either hERG inhibitors or non-inhibitors.



**SCHEME 1.** Schematical representation of the virtual experiment procedures

### Number of added compounds

| Exp. | 0 | 2 | 4 | 6 | 10 |
|---|---|---|---|---|---|
| **Inhibitors** | 0.46 (0.12) | 0.66 (0.49) | 0.80 (0.66) | 0.87 (0.71) | 0.91 (0.71) |
| | 0.68 (0.15) | 0.84 (0.49) | 0.89 (0.64) | 0.90 (0.56) | 0.94 (0.71) |
| | 0.42 (0.13) | 0.64 (0.38) | 0.82 (0.58) | 0.88 (0.70) | 0.92 (0.72) |
| | 0.42 (0.11) | 0.67 (0.34) | 0.81 (0.58) | 0.82 (0.34) | 0.90 (0.57) |
| | 0.59 (0.23) | 0.73 (0.33) | 0.88 (0.46) | 0.89 (0.47) | 0.94 (0.63) |
| | 0.53 (0.21) | 0.66 (0.29) | 0.77 (0.31) | 0.81 (0.29) | 0.86 (0.49) |
| | 0.68 (0.21) | 0.73 (0.24) | 0.84 (0.36) | 0.91 (0.50) | 0.88 (0.45) |
| **Non-inhibitors** | 0.11 (0.12) | 0.18 (0.39) | 0.24 (0.63) | 0.19 (0.60) | 0.19 (0.60) |
| | 0.09 (0.11) | 0.12 (0.45) | 0.17 (0.63) | 0.11 (0.79) | 0.14 (0.62) |
| | 0.12 (0.14) | 0.18 (0.42) | 0.24 (0.59) | 0.16 (0.75) | 0.19 (0.61) |
| | 0.12 (0.20) | 0.17 (0.43) | 0.20 (0.57) | 0.17 (0.57) | 0.10 (0.70) |

- – Inconclusive prediction
- – Predicted Non-inhibitor
- – Predicted Inhibitor

**NOTE:** The coloring scheme takes into account both predicted probability and the Reliability Index values

**TABLE 1.** Model performance for Test set compounds after different numbers of similar molecules added to the library (numbers in parentheses report prediction Reliability Index values)

## GALAS MODEL APPLICATION ON PUBCHEM DATA – AN EXAMPLE SCENARIO WITH CYP3A4 INHIBITION

GALAS model for the prediction of CYP3A4 enzyme inhibition developed at ACD/Labs using a training set of ca. 900 compounds was used as a starting point of this investigation. A recently published PubChem collection [2] containing more than 11,000 individual compounds was chosen as a good representation of an actual 'in-house' project for the external validation of ACD/Labs CYP3A4 inhibition model. For demonstration, the available PubChem data sets (cleaned from salts, mixtures, *etc.*) were classified using different thresholds:

- CYP3A4 inhibition in general (IC$_{50}$ < 50 uM) – 8528 compounds
- Effective CYP3A4 inhibition (IC$_{50}$ < 10 uM) – 7696 compounds

The first threshold corresponds to the criteria used in classification of the training set data of the ACD/Labs CYP3A4 inhibition model. The second threshold was introduced primarily considering the fact that there is actually no objective definition of what is a CYP3A4 inhibitor and as a result different classification schemes might exist. Additionally, even with the consistent classification threshold, a simple fact that a certain company is using property measurement protocol that is different from the ones usually used to measure the publicly reported values of the same property can still result in inconsistent qualitative data. All of these factors introduce additional data variability which is one of the causes contributing to the reduction of prediction quality.

Both PubChem sets have been split in half with one part of the compounds intended for the gradual addition to the blank Self-training Library, whereas the second one reserved for model performance evaluation.

Increasing size of PubChem based Self-training Library gives a steady rise in the number of test set compounds falling within the Applicability Domain of the model (RI>0.3) and obtaining high quality predictions (RI>0.5), which are correctly classified as positive or negative in terms of the property in all but a few cases, as shown in the previous example with hERG inhibition.
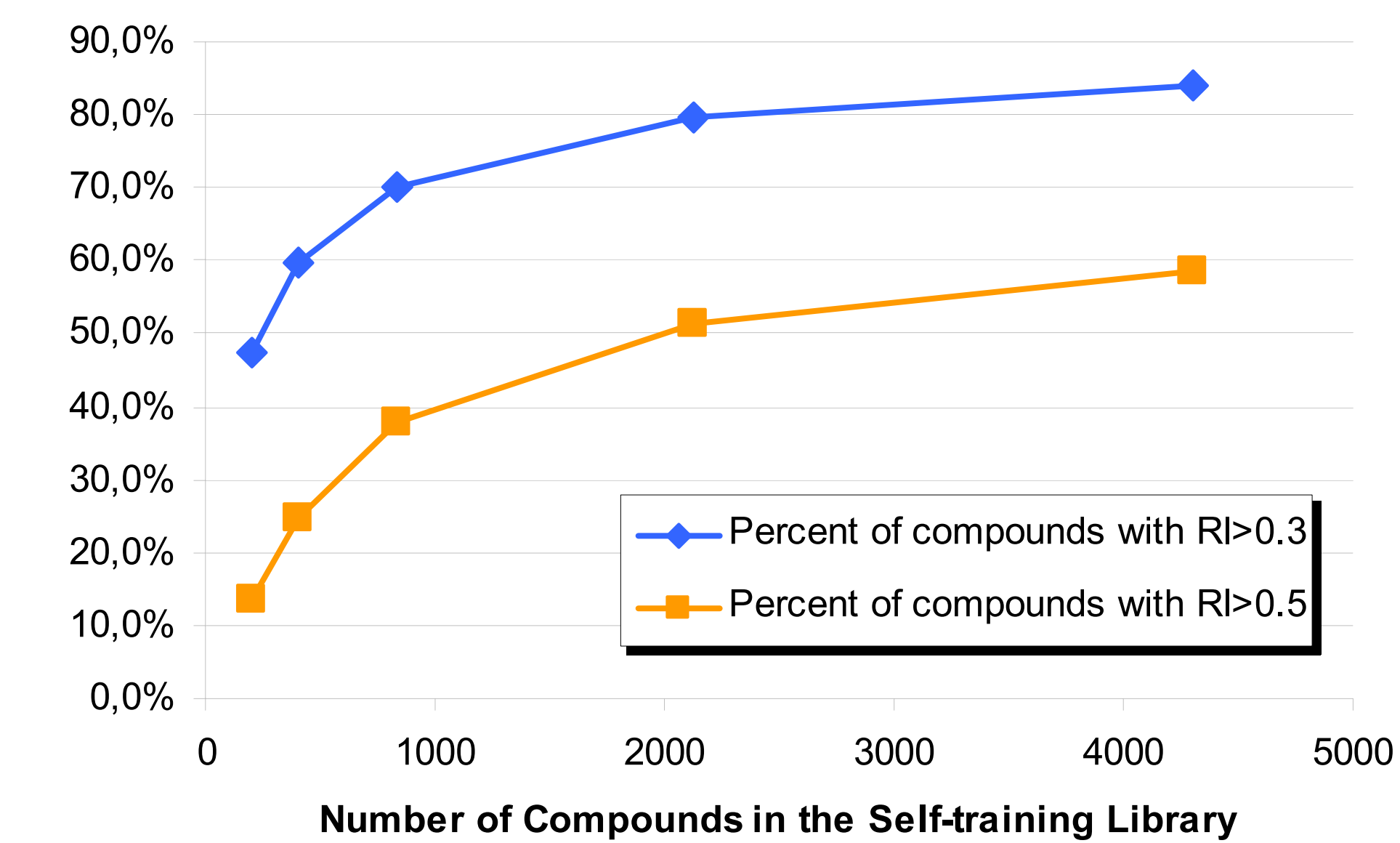


**FIGURE 2.** Number of corresponding reliability predictions following each addition of the general inhibition PubChem data to the Self-training Library of the CYP3A4 inhibition GALAS model
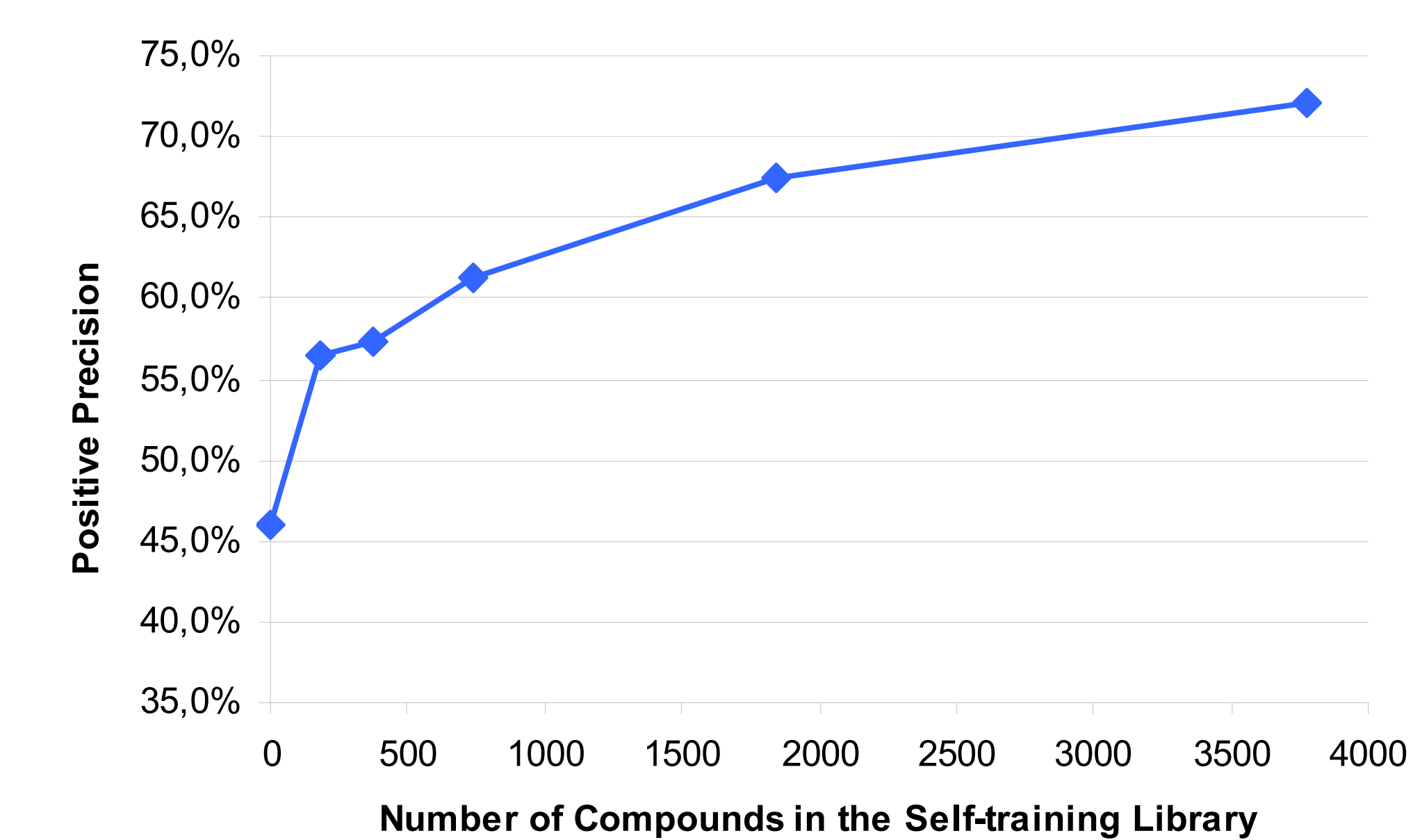
The positive precision (i.e., the fraction of true positives among all positive predictions of the model) of the initial ACD/Labs CYP3A4 inhibition model for the effective inhibition test set is ca. 40%. This is no surprise given the differences in classification thresholds used to obtain general and effective inhibition sets. However a dramatic impact on positive precision is observed if the first part of the effective inhibition set is used as a Self-training Library.

These observations suggest that the GALAS models can successfully cope with the practical challenges potentially arising during their applications in real life 'in-house' projects.



**FIGURE 3.** Changes in the positive precision of the GALAS model of CYP3A4 inhibition during its training with the effective inhibition PubChem set

## GALAS MODELS IN ADME AND TOX SUITES

Currently ACD/Labs software products contain trainable GALAS models for the following properties:

- Genotoxicity (Ames test)
- Acute rodent toxicity (LD$_{50}$)*
- Aquatic toxicity (LC$_{50}$)**
- P450 Substrate Specificity***
- P450 Inhibition Specificity***
- P-gp Substrate/Inhibitor Specificity
- hERG channel inhibition
- Plasma protein binding (LogKa and %PPB)
- Ionization constants (pKa)
- Quantitative solubility in pure water (LogS$_w$)
- Quantitative solubility in buffer (LogS)
- Qualitative solubility in buffer
- Octanol-water or buffer partitioning coefficients (LogP and LogD)

\* - Mouse OR, IP, IV, SC, and Rat OR, IP systems

\*\* - fathead minnow (Pimephales promelas), and water flea (Daphnia magna) species

\*\*\* - CYP3A4, CYP2D6, CYP2C9, CYP2C19, and CYP1A2 isoforms

## REFERENCES

[1] Zhu BY et al. *Bioorg Med Chem Lett.* **2006**, 16, 5507.

[2] *NCBI PubChem database*. Available at http://pubchem.ncbi.nlm.nih.gov/.