# NC STATE UNIVERSITY

# *Characterizing, Navigating, and Modeling the Chemical Space Using Next-Generation Cheminformatics Methods*
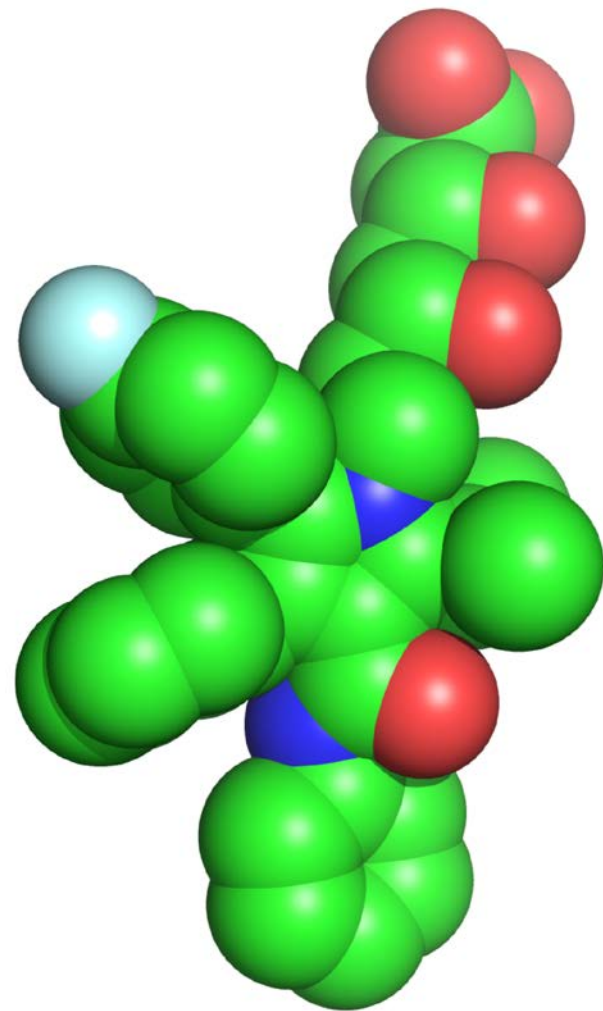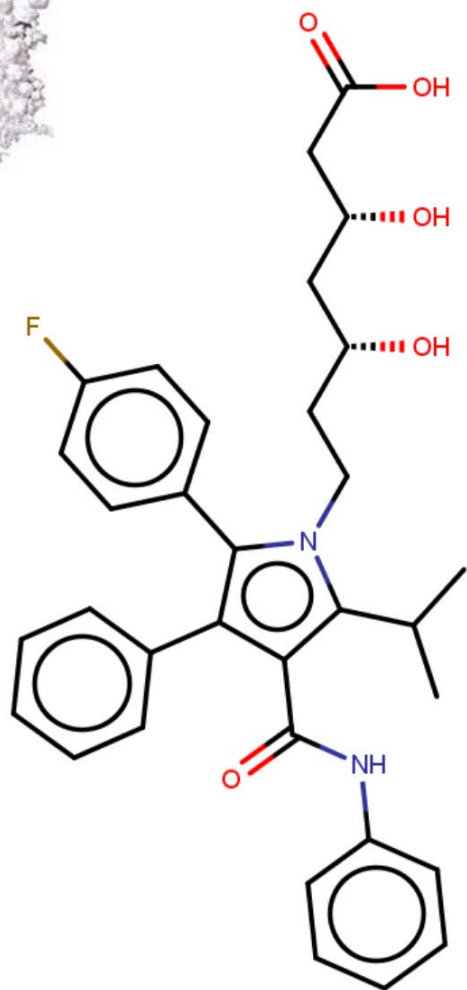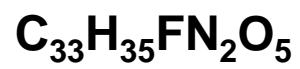
Denis Fourches, PhD

*Department of Chemistry, Bioinformatics Research Center,
North Carolina State University, USA
www.fourches-laboratory.com*
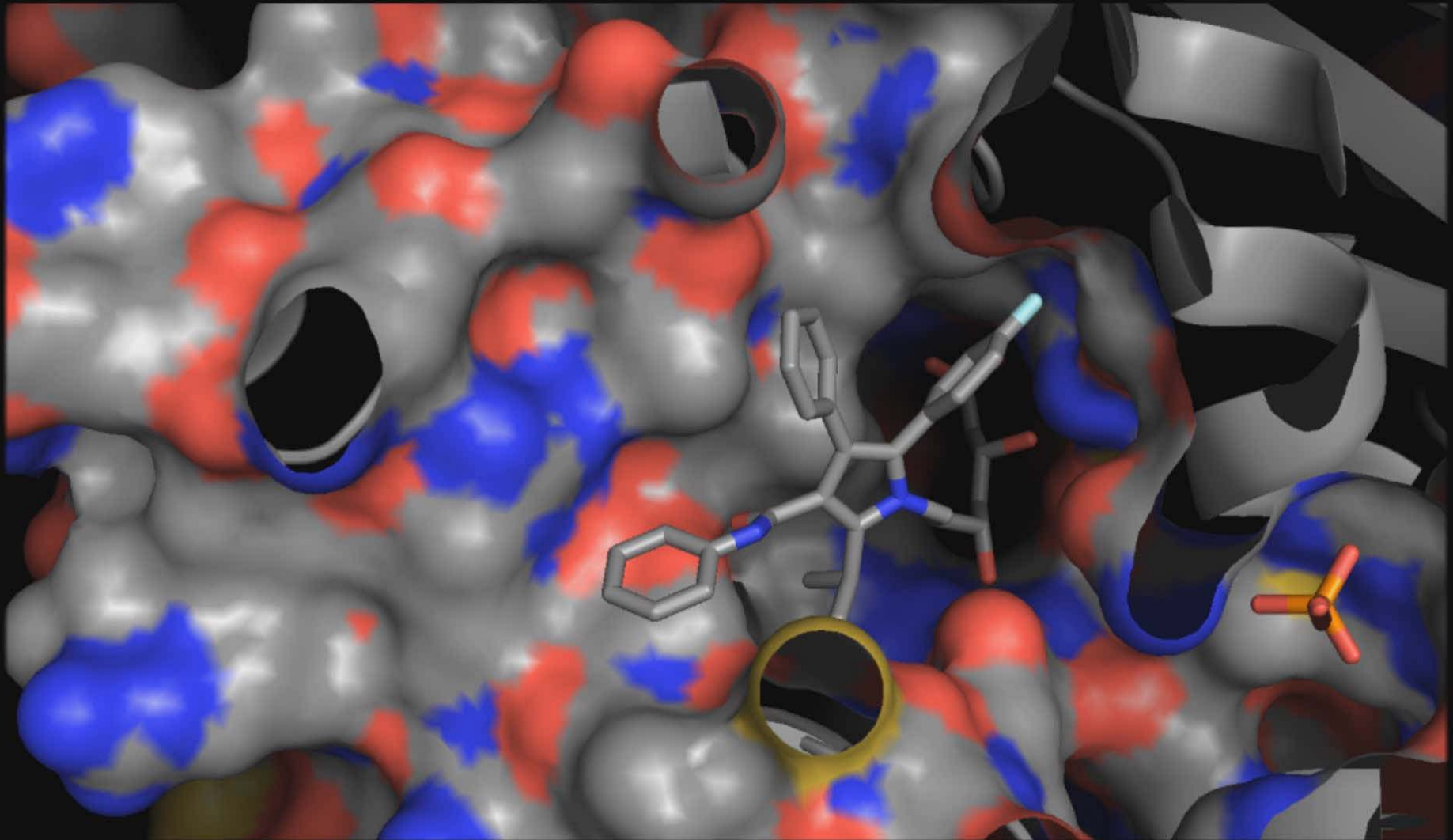
*Follow on*    ResearchGate    Linked in
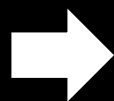
$C_{33}H_{35}FN_2O_5$

*0D/1D*

*2D*

*3D*

LIPITOR (atorvastatin)

World's best-selling drug of all time ($125 billion over 15 years)
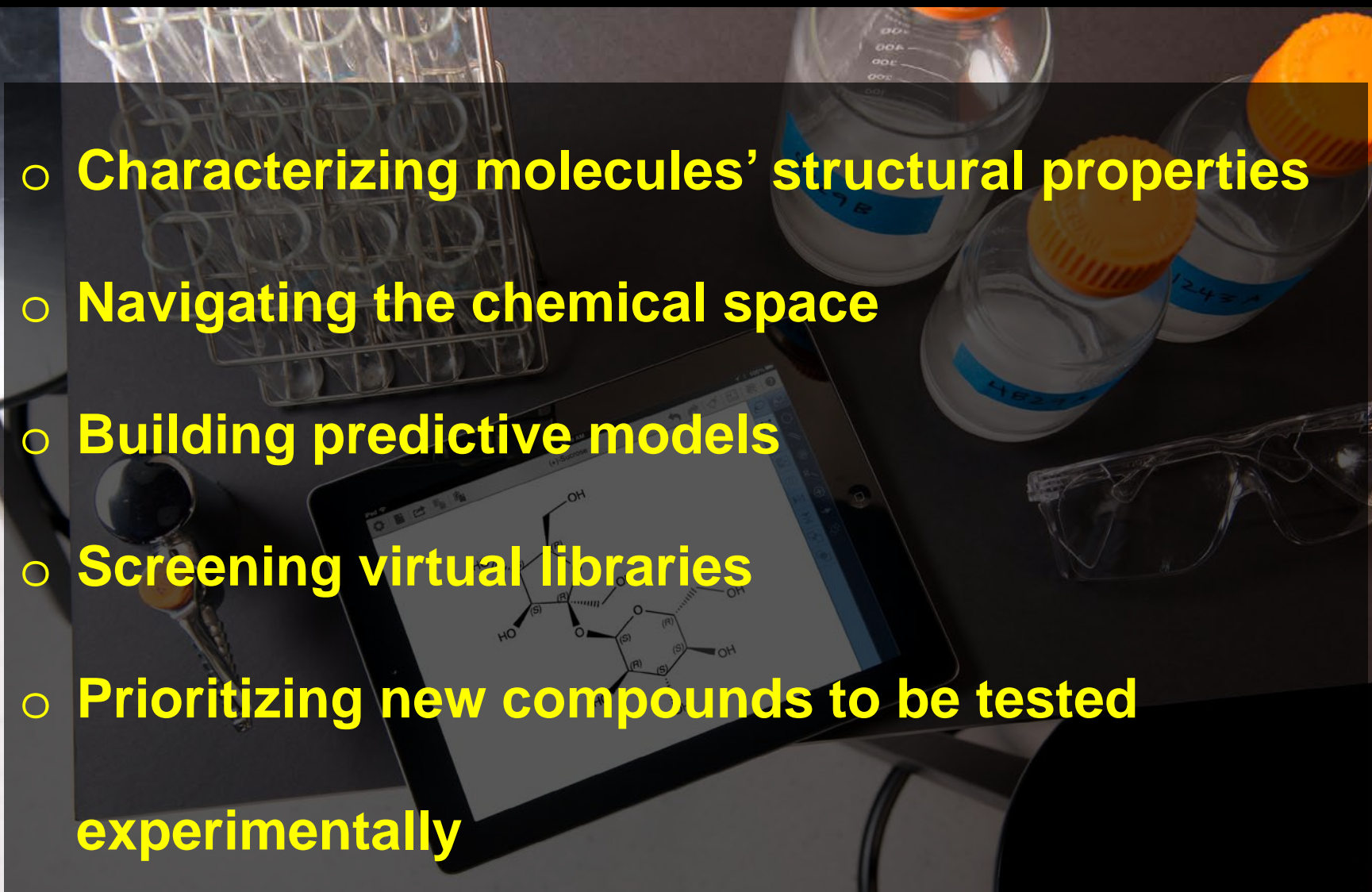
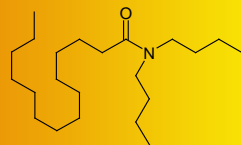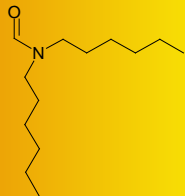Competitive inhibitor of HMG-CoA reductase in liver → Lower blood cholesterol

PDB code = 1HWK

# Cheminformatics is becoming an essential element in the chemist's toolbox

- Characterizing molecules' structural properties

- Navigating the chemical space

- Building predictive models

- Screening virtual libraries

- Prioritizing new compounds to be tested experimentally

**COMPOUNDS**

**DESCRIPTORS**

**ACTIVITY**

Thousands of molecular descriptors are available for organic compounds

constitutional, topological, structural, quantum mechanics based, fragmental, steric, pharmacophoric, geometrical, thermodynamical, conformational, etc.

0.613
0.380
-0.222
0.708
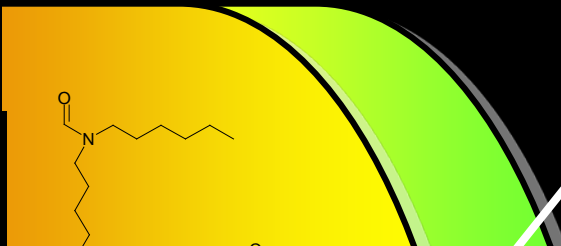1.146
0.491
0.301
0.141
0.956
0.256
0.799
1.195
1.005

- **Building of models** using machine learning methods (NN, SVM, RF)

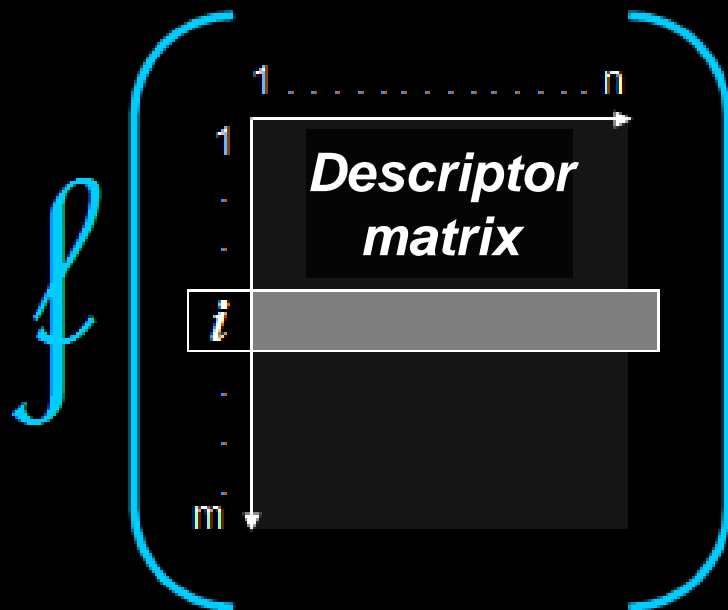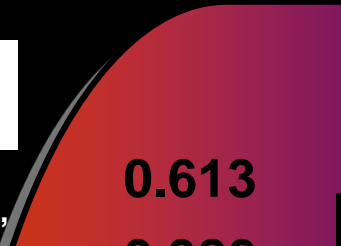- **Validation of models** according to numerous statistical procedures, and their applicability domains.

Cherkasov, Muratov, Fourches, et al. 2014, J Med Chem, 57(12), 4977-5010

**COMPOUNDS**

**ACTIVITY**

Thousands of molecular descriptors are available for organic compounds

constitutional, topological, structural, quantum mechanics based, fragmental, steric, pharmacophoric, geometrical

0.613

$$ f \begin{bmatrix} & 1 \dots \dots \dots n \\ 1 & \text{Descriptor matrix} \\ i & \\ m & \end{bmatrix} = \textbf{ACTIVITY} (i) $$

With m molecules and n descriptors

- **Validation of models** according to numerous statistical procedures, and their applicability domains.

1.005

Cherkasov, Muratov, Fourches, et al. 2014, J Med Chem, 57(12), 4977-5010

Correct Classification Rate (CCR) for QSAR models discriminating sensitizers from non-sensitizers was 71–88% when evaluated on several external validation sets, within a broad AD, with positive (for sensitizers) and negative (for non-sensitizers) predicted rates of 85% and 79% respectively.

Alves, Muratov, Fourches, Strickland, Kleinstreuer, Andrade, Tropsha. Toxicol Appl Pharmacol. 2015 Apr 15;284(2):262-72.

Skin penetrants

No penetration

QSAR analysis

Non-sensitizers

Sensitizers

# Curation of Chemogenomics Data

*Fourches, Muratov, Tropsha. Nature Chemical Biology. 2015, 11, 535.*

*Fourches, Muratov, Tropsha. JCIM. 2016, In Press.*

# Can any QSAR model, even if well validated, be applied to any molecule ?

| QSPR Models | ●──────▶ | Test compound |

## Prediction Performance

**Robustness of QSPR models**

- Descriptors type;
- Descriptors selection;
- Machine-learning methods;
- Validation of models.

**Applicability domain of models**

Is a *test compound* similar to the *training set* compounds?

# Applicability Domain of a given QSAR model

● = **TEST COMPOUND**

*Descriptor 2*

**TRAINING SET**

*Descriptor 1*

**INSIDE THE DOMAIN**

**OUTSIDE THE DOMAIN**

$AD_i \leq 100\ \%$

Will be predicted by the model

$AD_i > 100\ \%$

Will not be predicted by the model

The new compound will be predicted by the model, only if :

$$D_i \leq <D_k> + Z \times s_k$$

with Z, an empirical parameter (0.5 by default)

Tropsha et al., J. Med. Chem, 2002, 45, p. 2811-2823

**AD** parameter of applicability domain

$$AD_i = D_i / ( <D_k> + Z \times s_k ) \times 100$$
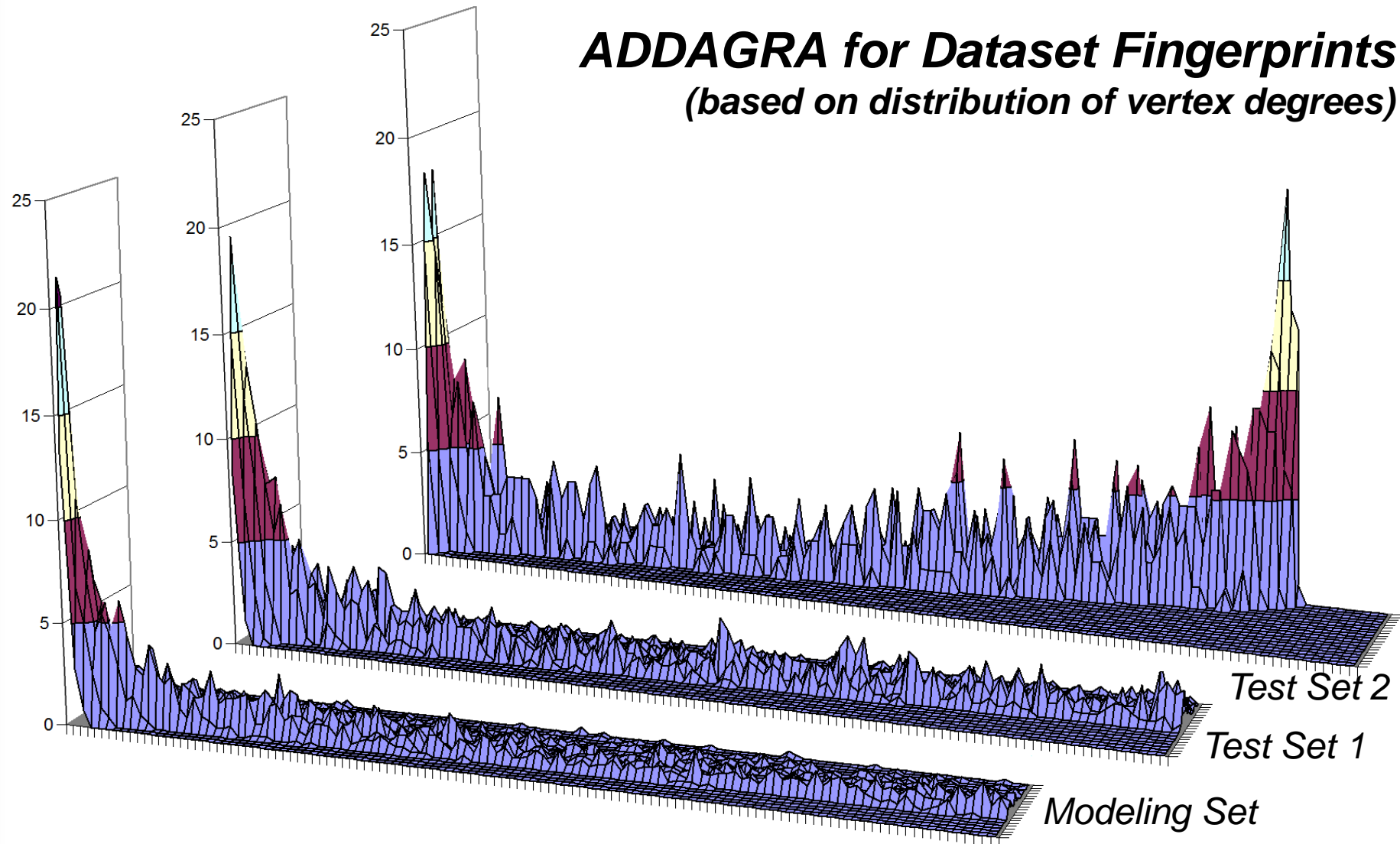
# Development of New Computational Tools Adapted to Hyper-Dimensional HTS data to Analyze Drugs' Polypharmacology



| | | | |
|---|---|---|---|
| 1 | OT | 9 | M1 |
| 2 | V2 | 10 | M3 |
| 3 | V1b | 11 | M3D |
| 4 | V1a | 12 | M5 |
| 5 | PAR1 | 13 | H1 |
| 6 | 5-HT2A | 14 | P2Y1 |
| 7 | 5-HT2B | 15 | P2Y2 |
| 8 | 5-HT2C | 16 | P2Y4 |
| | | 17 | P2Y6 |

| | |
|---|---|
| 18 | P2Y11 |
| 19 | NK1 |
| 20 | NK2 |
| 21 | NK3 |
| 22 | Alpha 1A |
| 23 | Alpha 1D |
| 24 | CCK2 |

# Visualizing and comparing chemical datasets using the ADDAGRA approach



*ADDAGRA for Dataset Fingerprints*
*(based on distribution of vertex degrees)*

Test Set 2

Test Set 1

Modeling Set

**D.Fourches** and A.Tropsha. Using Graph Indices for the Analysis and Comparison of Chemical Datasets. Molecular Informatics, 2013, 32, 827–842.

# **ADDAGRA** for the analysis of prediction outliers



Mol 192

Mol 492

Mol 413

Mol 504

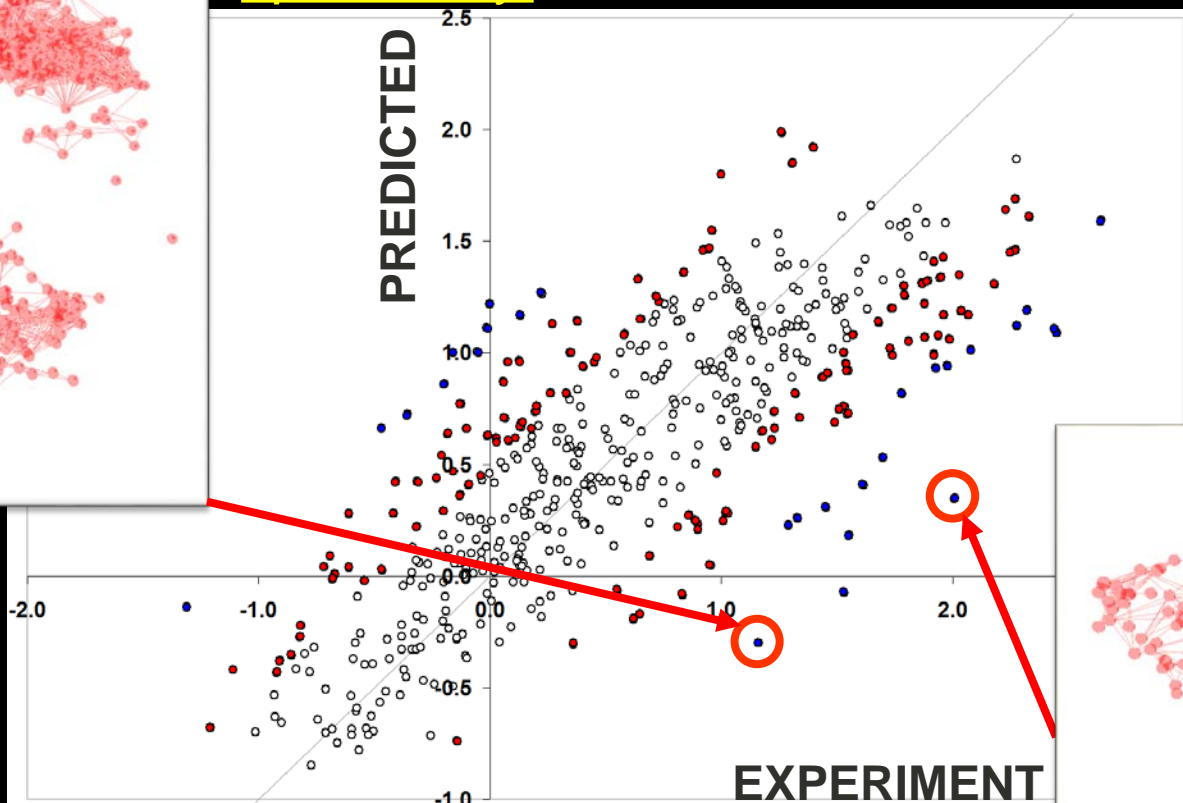**Many prediction outliers also correspond to outliers in the descriptor space.**

# **ADDAGRA** for the analysis of prediction outliers



**ACTIVITY CLIFF**

Mol 60

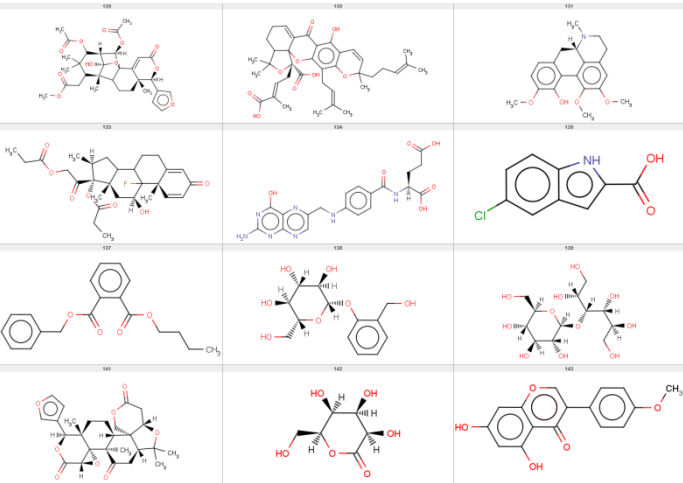But still, some large outliers cannot be identified in the chemical space only.

**MISANNOTATED**

Mol 441

PREDICTED

EXPERIMENT

# Hybrid modeling using both chemical and biological descriptors

**High Throughput Screening**

**BioAssay Summary**

Molecular weight, compositions and geometrical parameters, physico-chemical properties (acidic, basic, neutral, amphi- or lipophilic etc.)

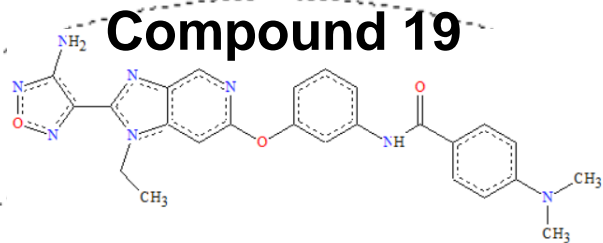**Molecular properties**

**Human health risk**

**Toxicity testing**

Rat

Mouse

**BIOLOGICAL DESCRIPTORS**

**CHEMICAL DESCRIPTORS**

*In Silico models*

16

**Compound 19**

**Chemical Neighbors**

*Structural Descriptors*

148

39

155

154

0.9

153

0.8

0.7

0.6

0.5

**CBRA Radial Plots**

*Low et al. CRT. 2013, 26(8):1199*

*Fourches et al. JCIM, 2016, In Preparation*

# Compound 19



**Biological Neighbors**
*GPCR BIOPROFILES*

**Chemical Neighbors**
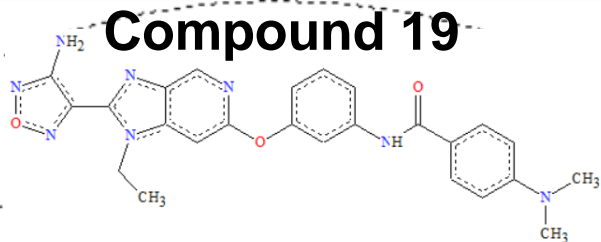*Structural Descriptors*

**Integrative Chemical Biological Read Across CBRA**

**CBRA Radial Plots**
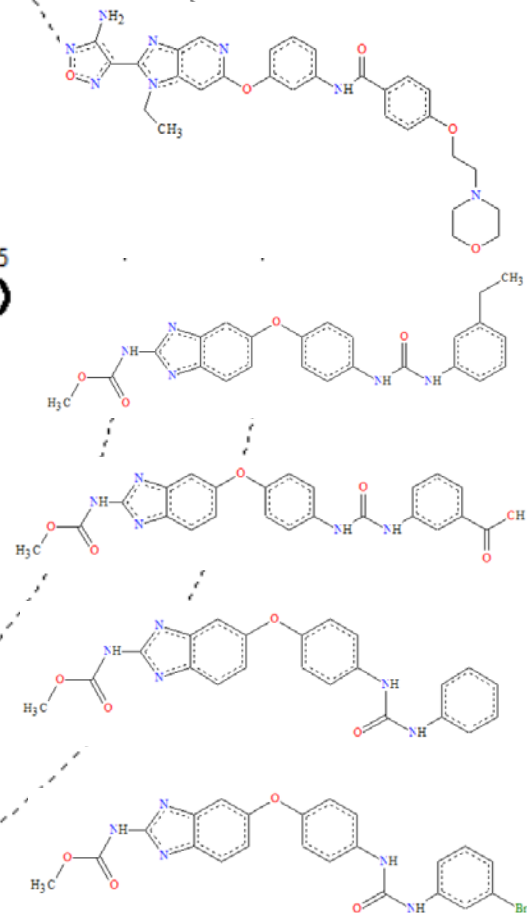*Low et al. CRT. 2013, 26(8):1199*

*Fourches et al. JCIM, 2016, In Preparation*

# Cheminformatics Approaches To Analyze the Similarity and the Effects of Environmental Chemical Mixtures

# Cheminformatics Approaches To Analyze the Similarity and the Effects of Environmental Chemical Mixtures

**CBRA Radial Plots Based on Both Chemical and Children's Exposures Similarity**

# Concept of Quantitative Structure-Exposure-Toxicity Relationships (QSETR)



**Need to develop new modeling workflow …**

# Molecular Docking of ERK2 Inhibitors

# GPU-accelerated molecular dynamics simulations

*Example: 1 ns simulation of ERK2 using Desmond*
*65Å*90Å*70Å, 42k atoms, explicit solvent (TIP3P water), step= 1 femtosecond*

## Up to 1 µs per day on high-end GPU workstations!

# New MD-QSAR modeling approach



~10³ Compounds

~10³

3D Descriptors

Time Steps

~10⁶

MD-based Descriptor Matrix

*Deep Learning* → **Hyper-Predictive QSAR models**

# Computation of MD Descriptors

- 3D descriptors computed for all conformations sampled along ligand MD simulations

- MD descriptors were constructed by taking the mean and standard deviation of each 3D descriptor distribution for each ligand:

$$\bar{x}_i = \frac{\sum_{j=1}^{n} x_{ij}}{n}$$

$$s_i = \sqrt{\frac{\sum_{j=1}^{n}(x_{ij} - \bar{x})^2}{n-1}}$$

Dataset: 87 ERK2 kinase inhibitors; pKi ranging from 4.6 to 9



**Atomic Masses Weighted WV**

Cohen's D = 1.1

Mean of 3D descriptor distribution

Cohen's D = .93

# Descriptor Set Distributions



- Compounds were classified as
  - Active : pKi >= 7.5
  - Inactive : pKi < 7.5

- For the 3D and MD descriptors, clear difference in the profiles of active and inactive compounds.

- For MACCS and 2D descriptors, the difference in active and inactive profiles is less apparent

**MACCS Fingerprints**

**MD Descriptors**

Cophenetic correlation coefficient: **0.89**

Cophenetic correlation coefficient: **0.74**

NC STATE UNIVERSITY

# Characterizing the MD Chemical Space
## of ERK2 Inhibitor Conformations

ACS Publications
MOST TRUSTED. MOST CITED. MOST READ.

ACS NANO

Fourches et al.
2010, 4(10): 5703-12.

Computational Descriptors

DESCRIPTORS

**Q**uantitative
**N**anostructure
**A**ctivity
**R**elationships

Nanoparticles 1 2 3 4

Activity Profiles

High-throughput
cellular-based assays

Experimental properties

NANOPARTICLES

Fourches D, Pu D, Tropsha A. *Comb Chem High Throughput Screen*. 2011
Fourches et al. *Nanotoxicology*, 2015, In Press.

# Functionalized Carbon Nanotubes



*Fourches et al.*
*Nanotoxicology. In Preparation*

# QNAR Modeling of Carbon Nanotubes

In 2008, Zhou et al* published *in vitro* protein binding, acute toxicity and immune toxicity assays for 84 Carbon NanoTubes (CNTs) decorated with different surface modifiers.



Different surface modifiers were introduced at the $R_1$, $R_1'$ and $R_2$ position

84 CNTs Tested in Two Different Types of Assays

Protein Binding

Toxicity

Bovine Serum Albumin
BSA

Carbonic Anhydrase
CA

Chymotripsin
CT

Hemoglobin
HB

Acute Toxicity

Immune Toxicity

*Zhou et al. Nano Lett., Vol. 8, No. 3, 2008*

# Computer-aided design of novel carbon nanotubes with desired biological properties
## (in collaboration with Dr. Bing Yan, St. Jude Children's Research Hospital)

~240,000 small molecules considered attachable to CNTs

QNAR models
Similarity filters
Empirical rules

VIRTUAL SCREENING

~10$^2$ molecules

Synthesis and experimental tests

A    B    C

100 μm

D    Normal          II-8 CNT          II-11 CNT

| CNT ID | II-1 | II-2 | II-3 | II-4 | II-5 | II-6 | II-7 | II-8 | II-9 | II-10 |
|---|---|---|---|---|---|---|---|---|---|---|
| Average cell viability (%) | 58 | 61 | 61 | 56 | 58 | 65 | 68 | 72 | 68 | 59 |
| Standard Deviation (%) | 5 | 3 | 3 | 3 | 2 | 10 | 6 | 7 | 3 | 6 |
| Experiment | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Predicted | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| CNT ID | II-11 | II-12 | II-13 | II-14 | II-15 | II-16 | II-17 | II-18 | II-19 | II-20 |
| Average cell viability (%) | 39 | 49 | 46 | 49 | 52 | 41 | 51 | 49 | 55 | 50 |
| Standard Deviation (%) | 9 | 8 | 7 | 5 | 8 | 11 | 5 | 9 | 11 | 10 |
| Experiment | 1 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 0 |
| Predicted | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |

| CNT ID | II-21 | II-22 | II-23 | II-24 | II-25 | II-26 | II-27 | II-28 | II-29 | II-30 |
|---|---|---|---|---|---|---|---|---|---|---|
| Average protein binding (F0/F1) | 1.77 | 1.78 | 1.87 | 1.76 | 1.82 | 2.30 | 1.74 | 2.00 | 1.67 | 2.33 |
| Standard Deviation | 0.05 | 0.06 | 0.02 | 0.03 | 0.02 | 0.02 | 0.02 | 0.01 | 0.06 | 0.02 |
| Experiment | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 |
| Predicted | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| CNT ID | II-31 | II-32 | II-33 | II-34 | II-35 | II-36 | II-37 | II-38 | II-39 | II-9 |
| Average protein binding (F0/F1) | 3.40 | 2.28 | 2.04 | 2.22 | 2.17 | 2.95 | 2.08 | 2.25 | 2.65 | 2.24 |
| Standard Deviation | 0.03 | 0.05 | 0.08 | 0.01 | 0.04 | 0.00 | 0.02 | 0.06 | 0.05 | 0.11 |
| Experiment | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| Predicted | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |

# Computer-aided design of novel carbon nanotubes with desired biological properties
(in collaboration with Dr. Bing Yan, St. Jude Children's Research Hospital)



**All** rationally prioritized, synthesized, and tested CNTs predicted as non-toxic were confirmed experimentally.

**6 out of 10** rationally prioritized, synthesized, and tested CNTs predicted as toxic were confirmed experimentally

QNAR models

Synthesis and experimental tests

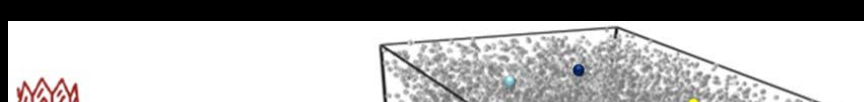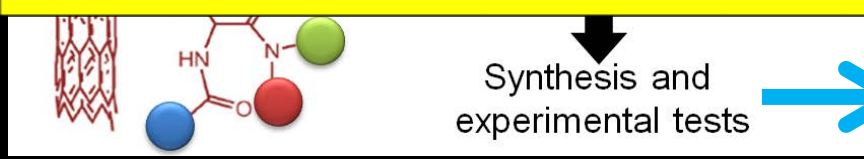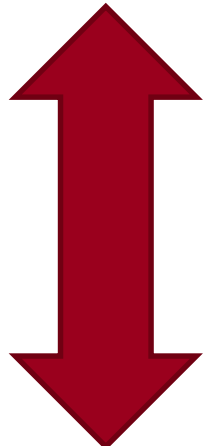| Standard Deviation (%) | 5 | 3 | 3 | 3 | 2 | 10 | 6 | 7 | 3 | 6 |
|---|---|---|---|---|---|---|---|---|---|---|
| Experiment | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Predicted | | | | | | | | | | |

| Experiment | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|
| Predicted | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **CNT ID** | II-31 | II-32 | II-33 | II-34 | II-35 | II-36 | II-37 | II-38 | II-39 | II-9 |
| Average protein binding (F0/F1) | 3.40 | 2.28 | 2.04 | 2.22 | 2.17 | 2.95 | 2.08 | 2.25 | 2.65 | 2.24 |
| Standard Deviation | 0.03 | 0.05 | 0.08 | 0.01 | 0.04 | 0.00 | 0.02 | 0.06 | 0.05 | 0.11 |
| Experiment | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| Predicted | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |

Fourches et al. Nanotoxicology 2016 In Press.
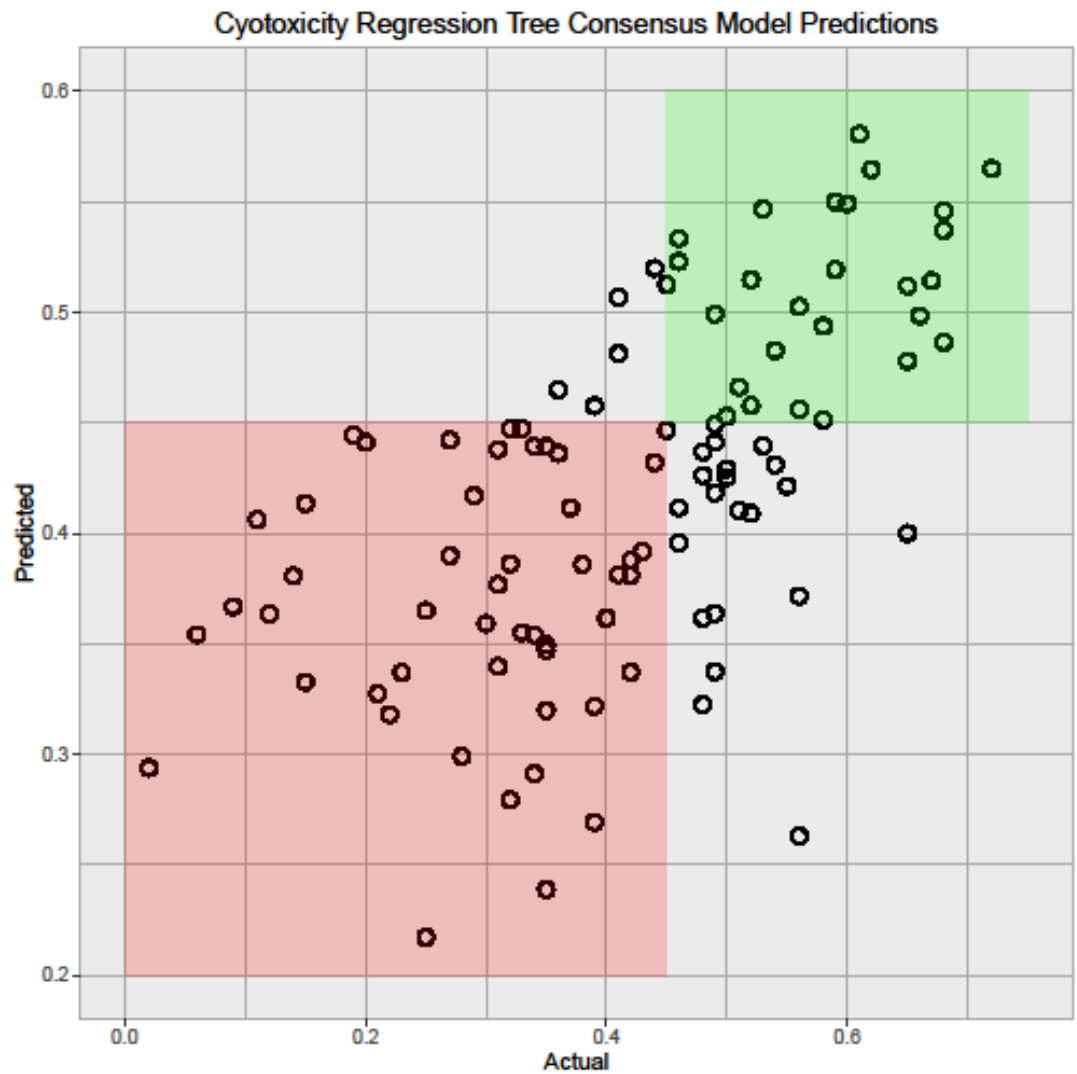
Fourches et al.
Nanotoxicology. In Preparation

# QNAR models for the enhanced set of f-CNTs

**Classification Models (LOO)**

| Cytotoxicity | Threshold | Sensitivity | | Specificity | | Accuracy |
|---|---|---|---|---|---|---|
| Random Forest | 0.45 | 79.25% | 42 | 62.00% | 31 | 70.87% |
| SVM-C Model | 0.45 | 74.40% | 41 | 56.00% | 28 | 66.99% |

**Continuous Models (LOO)**



Cyotoxicity Regression Tree Consensus Model Predictions
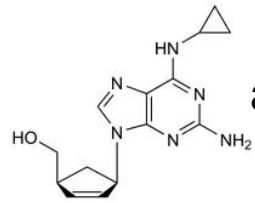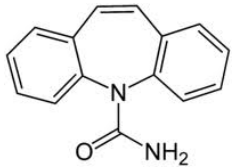
# HLA-induced drug adverse effects

**Small molecule drugs bind specifically to certain type of HLA proteins.**

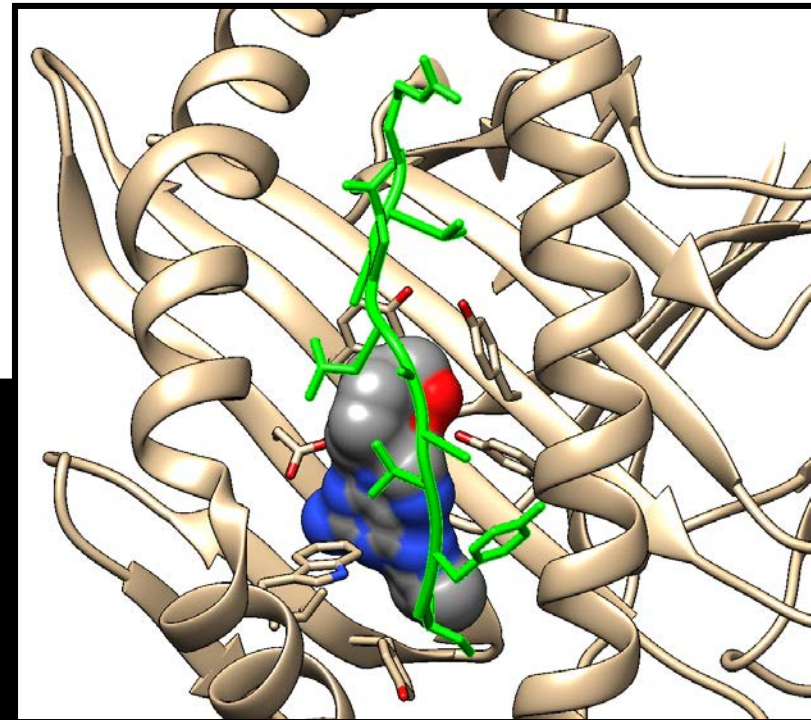abacavir → HLA-B*57:01 → Abacavir Hypersensitivity Syndrome (AHS)

carbamazepine → HLA-B*15:02 → Stevens-Johnson Syndrome (SJS)

➤ Binding pocket of HLA proteins is slightly modified

➤ Range of "self" peptides able to bind the given HLA type is modified
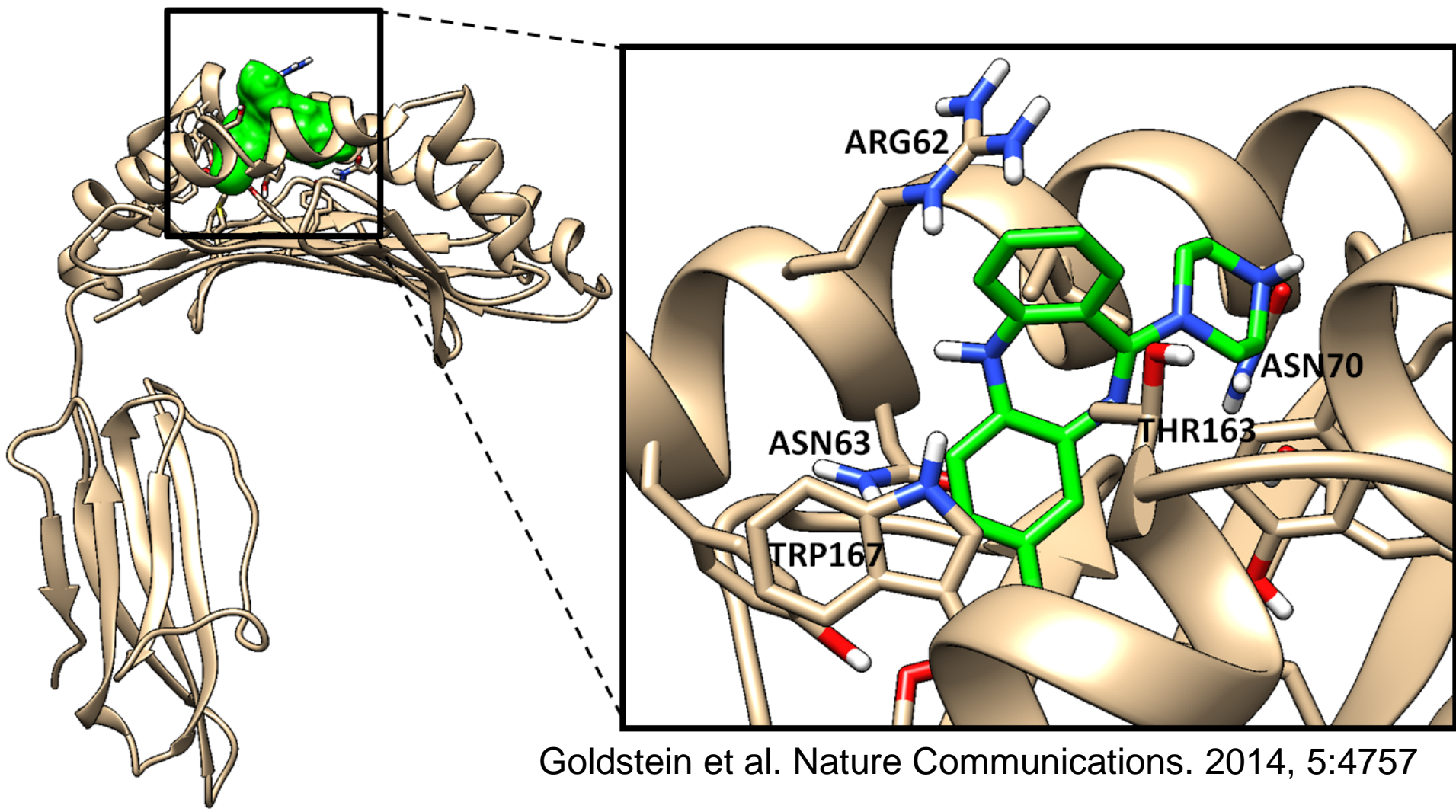
➤ Immune reaction

➤ **HLA proteins as important off-targets**

➤ **Need predictive models to assess HLA-induced ADR**

**Clozapine-induced agranulocytosis is associated with rare HLA-DQB1 and HLA-B alleles.**
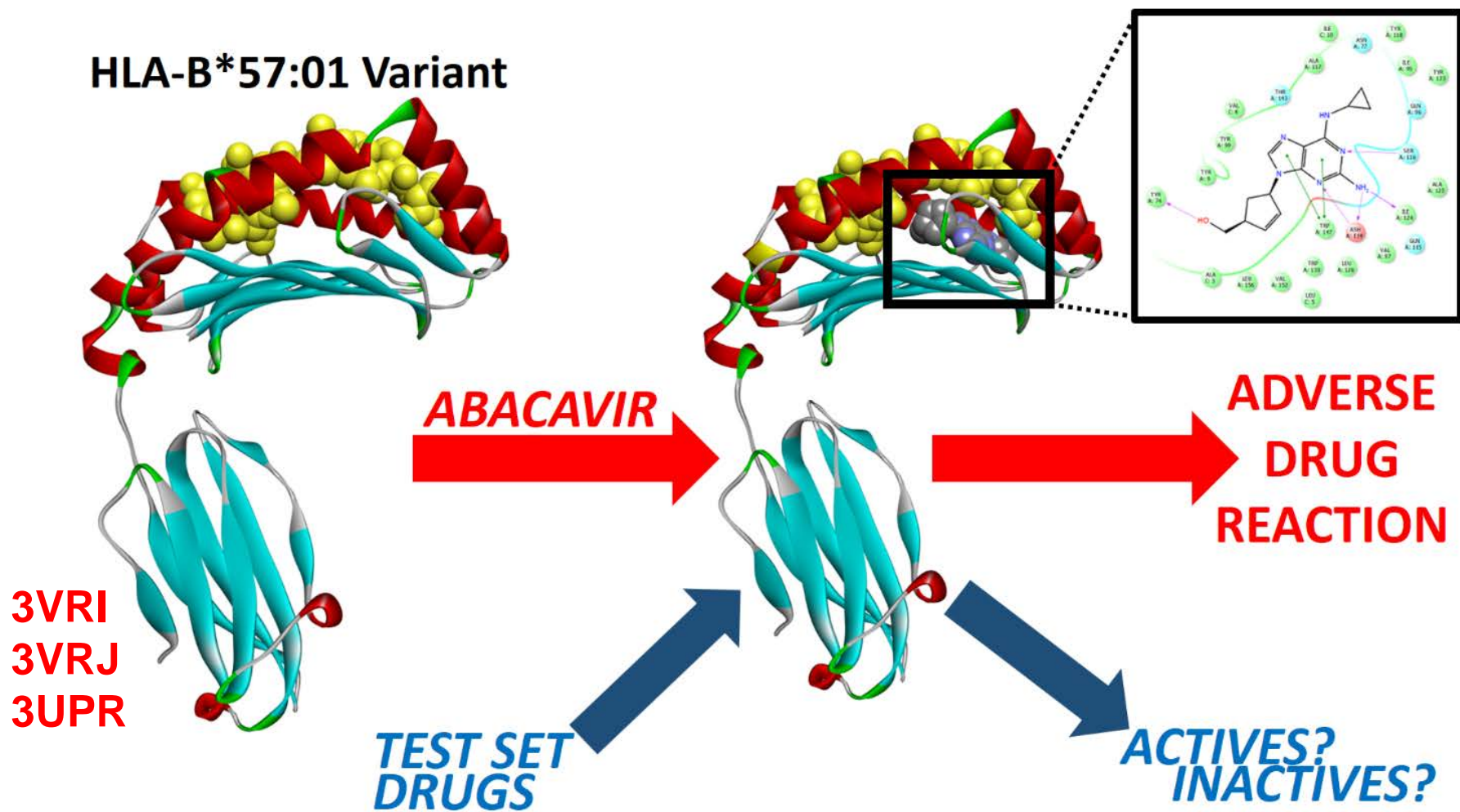
**GWAS + Docking**

Goldstein et al. Nature Communications. 2014, 5:4757
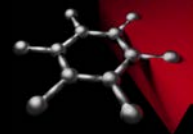
Homology model of HLA-B39

# Molecular Docking Study at HLA-B*57:01

Develop a virtual screening model using molecular docking at the HLA-B*57:01 variant using three X-ray crystals of abacavir bound to HLA-B*57:01.

# Summary

- Cheminformatics is becoming mandatory in all projects involving the curation, integration, characterization, analysis, testing, modeling, visualization, screening of chemicals.

- The skyrocketing amount of freely-available data in the public domain is boosting the development of new cheminformatics approaches to fully exploit that data, especially when it comes to chemical risk assessment and *in silico* toxicity predictions.

- New methods such as MD-QSAR or QSETR are poised to boost the prediction performances of cheminformatics predictors.

- Structure-based docking and pan-target screening have never been so relevant for chemical risk assessment, especially for key targets such as ER, AhR, or HLA.

# Acknowledgements



**Lab members**

George Van Den Driessche

Jeremy Ash

Melaine Kuenemann, PhD

Ryan Lougee

Phyo Phyo Kyaw Zin

Bethany Cook