# Applications of machine learning and artificial intelligence to designing chemicals and materials with the desired properties
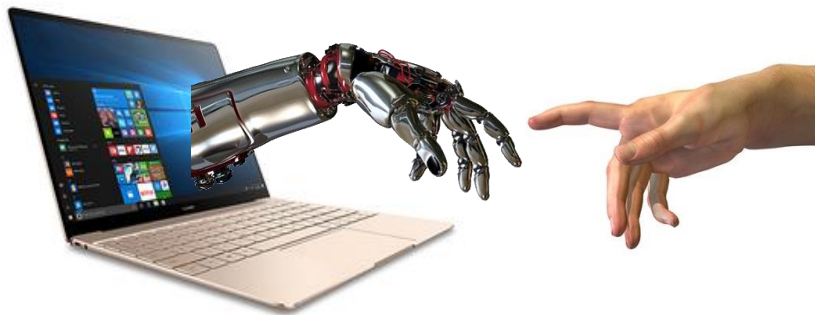
**Alexander Tropsha**

UNC Eshelman School of Pharmacy

# Outline

- Brief notes on machine learning/QSAR
- Materials Informatics and Materials Design
- Design, development and application of the <u>Re</u>inforcement <u>Lea</u>rning for <u>S</u>tructural <u>E</u>volution (ReLeaSE)*
- Summary and future work: QSAR without borders

*Popova, Mariya, Olexandr Isayev, and Alexander Tropsha. "Deep reinforcement learning for de-novo drug design." *arXiv preprint arXiv:1711.10907* (2017); Science Advances, 2018, in press
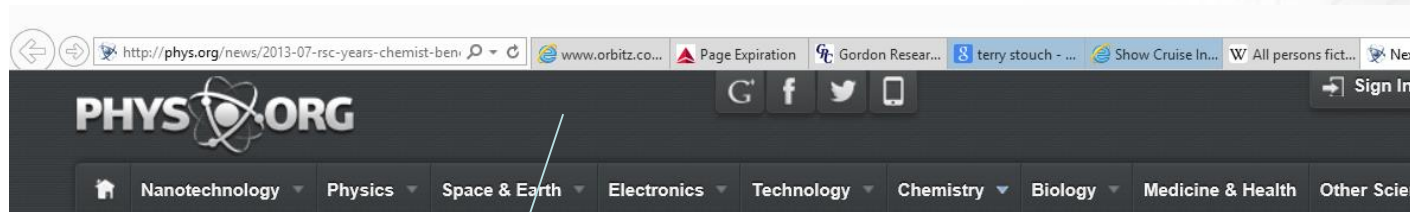
# Machine Learning Framework

$$y = f(\mathbf{x})$$

output    prediction function    Molecular features

- **Training:** given a *training set* of labeled examples $\{(\mathbf{x}_1, y_1), ..., (\mathbf{x}_N, y_N)\}$, estimate the prediction function $f$ by minimizing the prediction error on the training set

- **Testing:** apply $f$ to a never before seen *test example* $\mathbf{x}$ and output the predicted value $y = f(\mathbf{x})$

# The growing appreciation of molecular modeling and informatics

# Automated Retrosynthesis (Chematica)

# The growing appreciation of molecular modeling and informatics

# Promise of dramatic acceleration of drug discovery

**Pharma VOICE.com**

READ. THINK. PARTICIPATE.

| News | Blog | R&D | Commercial | Operations | |
|---|---|---|---|---|---|
| Magazine | PharmaVOICE 100 | | Resources | Events | Editorial | Advertise | Subscribe |

## GSK Has Developed A New Analytics Platform That Can Reduce The Time It Takes To Analyze Clinical Data From Months To Clicks

Source:
Thomas Macaulay, CIO UK
March 12, 2018

The platform uses large-scale data analytics to drive better decisions about the drug discovery pipeline, by allowing the pharmaceuticals giant to test the potential for new drugs before it begins clinical trials.

# Rise of the machines in legal industry



legal**it** insider
www.legaltechnology.com

## Deloitte Insight: Over 100,000 legal roles to be automated

*Added on the 16th Mar 2016 at 10:28 am*

Over 100,000 jobs in the legal sector have a high chance of being automated in the next twenty years, according to extensive new analysis by Deloitte.

The Deloitte Insight report, which predicts "profound reforms" across the legal profession within the next 10 years, finds that 39% of jobs (114,000) in the legal sector stand to be automated in the longer term as the profession feels the impact of more "radical changes."

# The ultimate dream of a computational chemist

# The chief utility of computational models: Annotation of new compounds



CHEMICAL STRUCTURES → CHEMICAL DESCRIPTORS → PREDICTIVE QSAR MODELS ← PROPERTY/ACTIVITY

QSAR MAGIC

CHEMICAL DATABASE

$\sim 10^6 - 10^9$ molecules

VIRTUAL SCREENING → Confirmed actives (toxic)

Confirmed inactives (non-toxic)

10

# QSAR Modeling Workflow: the importance of rigorous validation

**Datasets**

**Experimental confirmation**

Virtual screening (with **AD threshold**)

External set → Evaluation of external performance

An ensemble of QSAR Models

Modeling set

Internal validation Model selection

*courtesy of L. Zhang*

5-fold External Validation

1 5 2 3 4

*M o d e l i n g   m e t h o d s*

**Combi-QSAR modeling**

| *K*-Nearest Neighbors (*k*NN) | Random Forest (RF) | Support Vector Machines (SVM) |

*D e s c r i p t o r s*

Dragon    MOE

Tropsha, A. Best Practices for QSAR Model Development, Validation, and Exploitation *Mol. Inf.,* **2010**, 29, 476 – 488
*Fully implemented on CHEMBENCH.MML.UNC.EDU*

11

# Material Science and the Rise of Materials Informatics

- Explosive growth of materials data, both experimental databases and computational repositories.

  - <u>Structural data</u>: 160,000 entries in the Inorganic Crystal Structure Database (ICSD)

  - <u>Experimental data</u>: Numerous commercial and open experimental databases NIST, MatWeb, MatBase etc.

  - <u>Computational data</u>: Huge databases such as AFLOWLIB, Materials Project, and Harvard Clean Energy

  - Chemical space of possible materials is HUGE $\sim 10^{100}$ candidates [*Nat. Chem*. 7, 274-275 (**2015**)]

- Materials Genome Initiative or MGI (US Govt): Need for new high performance materials

# Closing the gap: materials structure-property relationships

**Materials**

**x –ray diffraction pattern**

**property**

**theoretical fingerprint**

**modeling**

predictive QSPR modeling

Quantitative Structure Activity Relationship approaches (QSAR)
Quantitative Structure-Property Relationships (QMSPR)

Isayev, Fourches, Muratov, Oses, Rasch, Tropsha, Curtarolo, Materials Cartography: Representing and Mining Materials Space Using Structural and Electronic Fingerprints. Chem. Mater., **2015**, 27: 735–743

# Material Informatics/MQSAR Workflow

Isayev, Fourches, Muratov, Oses, Rasch, Tropsha, Curtarolo, Materials Cartography: Representing and Mining Materials Space Using Structural and Electronic Fingerprints. Chem. Mater., **2015**, 27: 735–743

# Material Map (B-Fingerprints)

>15000 materials from ICSD
DFT PBE calculations from aflowlib.org



Orphans

Cluster C: metallic comp. with non-metallic atoms

Cluster B: bimetals, polymetals

Cluster D: small band gap comp., semiconductors

Orphans

Cluster A: insulators, ceramics, complex oxides

Band gap, eV

1          10

Isayev, Fourches, Muratov, Oses, Rasch, Tropsha, Curtarolo, Materials Cartography: Representing and Mining Materials Space Using Structural and Electronic Fingerprints. Chem. Mater., **2015**, 27: 735–743

# Systematic representation of materials using fragment descriptors



**A. Crystal Structure**

**B. Voronoi tessellation and neighbors search**

**C. Infinite periodic graph construction and property labeling (EA, IP, En, Rcov, etc)**

**D. Decomposition to fragments**

Nodes (atoms)

Edges (bonds, vDW contacts)

Path fragments of length N, N = 2, 3, …

Circular fragments (polyhedrons)

Isayev et al. Universal Fragment Descriptors for Predicting Properties of Inorganic Crystals. *Nature Comm*, **2017**, 8, 15679.

# ML Workflow for Materials Property Prediction



**Electronic Properties**

no $E_{BG}$

yes

**crystal structure**

classification model

metal or insulator?

no

band gap energy prediction

regression model

$\{E_{BG} \in \mathbf{R} : E_{BG} > 0\}$

Only!

**Thermo-Mechanical Properties**

bulk modulus (VRH) prediction

$\{X \in \mathbf{R}\}$

regression models

- Bulk modulus
- Shear modulus
- Thermal expansion
- Heat Capacity
- Thermal conductivity, etc.

All models are trained based on DFT-computed properties (VASP s/w from U. Vienna)

Isayev et al. Universal Fragment Descriptors for Predicting Properties of Inorganic Crystals. *Nature Comm*, **2017**, 8, 15679.
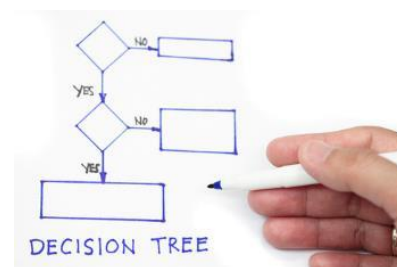
# Prediction of Electronic Properties

## electronic properties

a)



b)



Classification accuracy 95%
ROC Curve (AUC) 0.98

**Learning approach for all models:**
Gradient Boosting Decision Trees (GBT)

Isayev et al. Universal Fragment Descriptors for Predicting Properties of Inorganic Crystals. *Nature Comm*, **2017**, 8, 15679.

# Prediction of Thermomechanical Properties

($E_{BG}$ - band gap energy; $B_{VRH}$ - bulk modulus; $G_{VRH}$ -shear modulus; $t_D$ - Debye temperature; $C_P$ - heat capacity at constant pressure; $C_V$ - heat capacity at constant volume; $a_V$ -thermal expansion coefficient

| property | RMSE | MAE | $r^2$ |
|---|---|---|---|
| $E_{BG}$ | 0.51 (eV) | 0.35 (eV) | 0.90 |
| $B_{VRH}$ | 14.25 (GPa) | 8.68 (GPa) | 0.97 |
| $G_{VRH}$ | 18.43 (GPa) | 10.62 (GPa) | 0.88 |
| $\theta_D$ | 56.97 (K) | 35.86 (K) | 0.95 |
| $C_P$ | 2.31 ($k_B$/cell) | 0.84 ($k_B$/cell) | 0.99 |
| $C_V$ | 2.01 ($k_B$/cell) | 0.70 ($k_B$/cell) | 0.99 |
| $\alpha_V$ | $1.47 \times 10^{-5}$ (K)$^{-1}$ | $5.69 \times 10^{-6}$ (K)$^{-1}$ | 0.91 |

TABLE I. Statistical summary of the *five-fold cross-validated predictions* for the seven regression models (Figure 3).

Isayev et al. Universal Fragment Descriptors for Predicting Properties of Inorganic Crystals. *Nature Comm*, **2017**, 8, 15679.
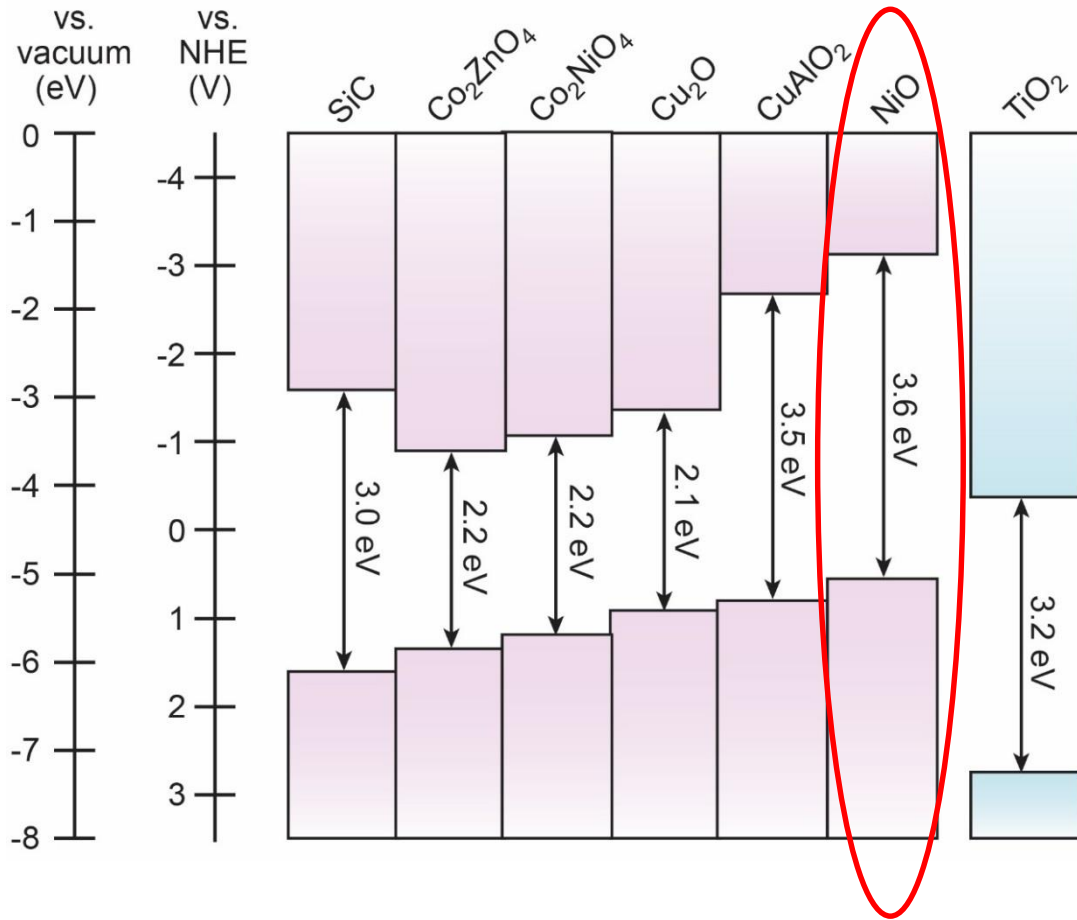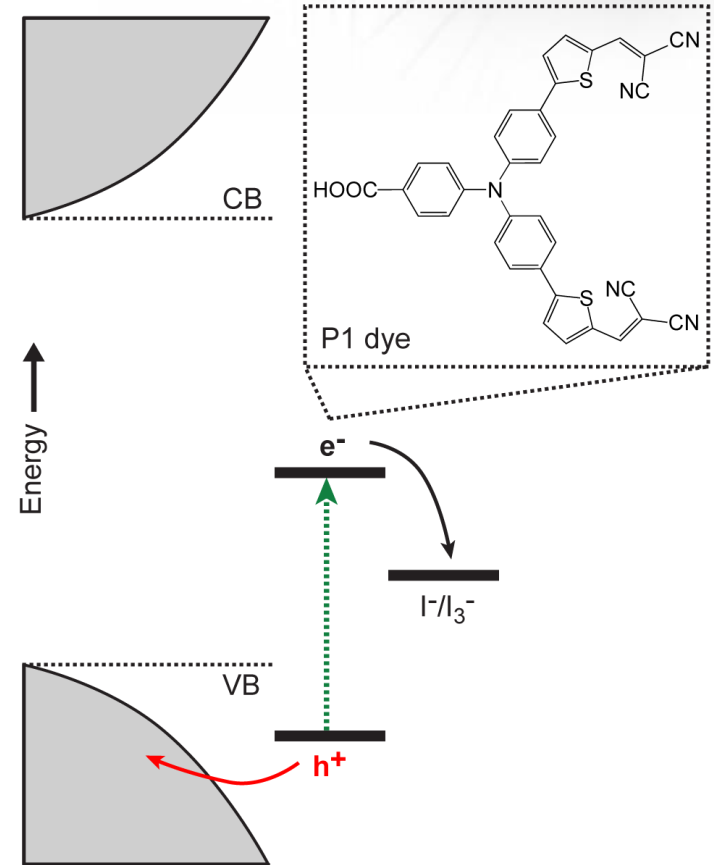
# Summary of Materials Informatics: Methods

- Fast, accurate general purpose machine learning methods for material's property prediction. Milliseconds on laptop vs. days on HPC cluster

- Universal applicability to different materials: currently covered 85 elements (H – Pu, without noble gases, Tc, Fr, Ra). All types of crystal lattices and symmetries.
  - Most competing approaches are specific to one prototype/family of materials or single property

- Works for other properties: elastic, thermoelectric, etc.

- Possible to gain *some* chemically/physically interpretable insight into "black box" model.

- Possible to derive materials design rules

- User friendly web app and RESTful API (http://aflow.org/aflow-ml/)

# Photocathode materials
## Evaluated as DSSCs
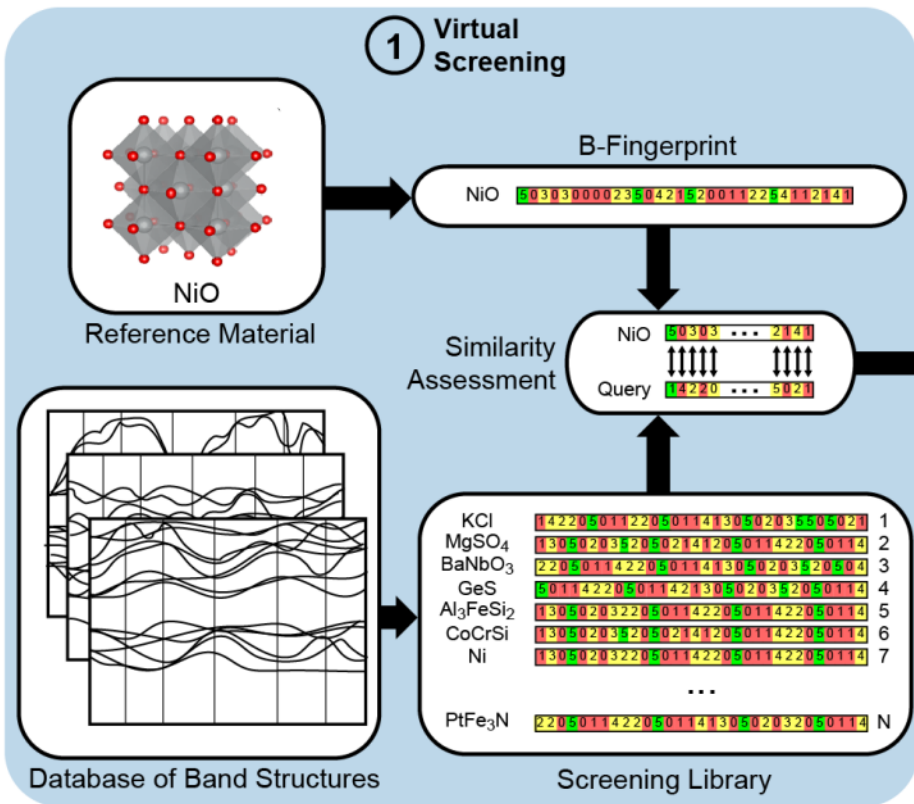
Dye-sensitized solar cells (DSSCs)

# Design of alternate photocathodes
## A materials informatics approach



(AFLOWLIB)

Moot, Isayev, Tropsha, Cahoon, Materials Discovery, **2016**, 6, 9-16

# Design of alternate photocathodes
## A materials informatics approach



(AFLOWLIB)

$PbTiO_3$ was identified as very similar to NiO
AND
**It is has a dielectric constant >100**

Moot, Isayev, Tropsha, Cahoon, Materials Discovery, **2016**, 6, 9-16

# Design of alternate photocathodes
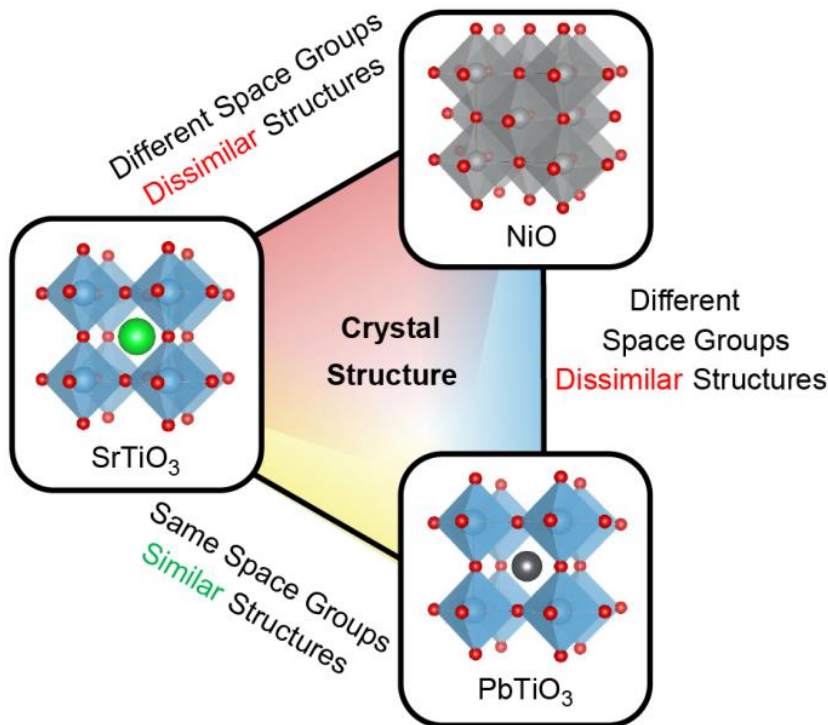## A materials informatics approach



(AFLOWLIB)

PbTiO$_3$ was identified as very similar to NiO
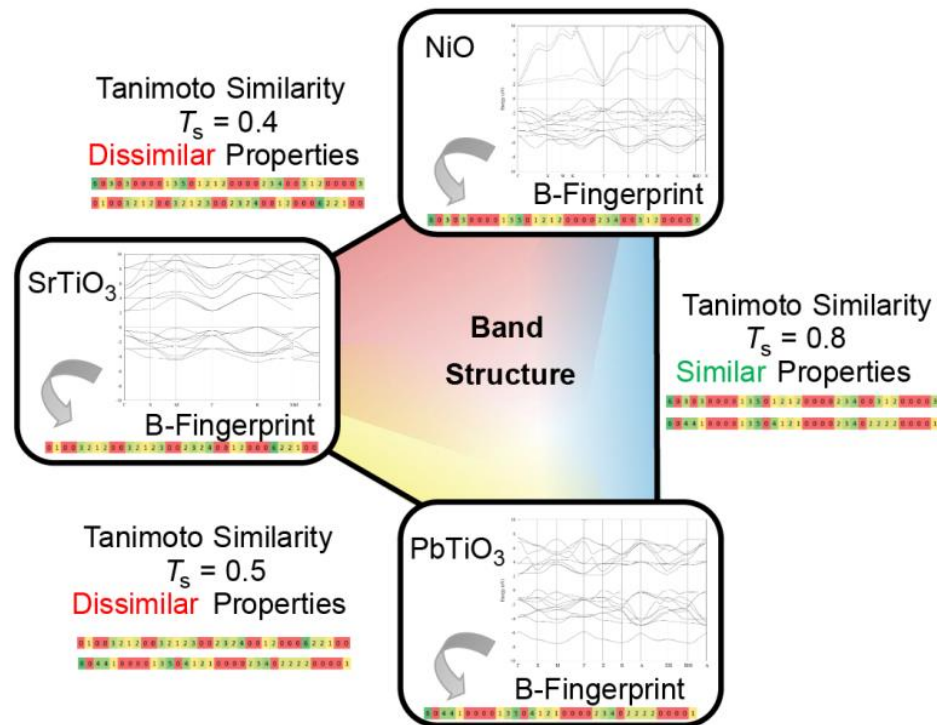AND
**It is has a dielectric constant >100**

# Materials informatics
## Identifying top hit: PbTiO$_3$

**Structure**

**Properties**



PbTiO$_3$ was identified as very similar to NiO in terms of electronic properties despite different crystal structures

Moot, Isayev, Tropsha, Cahoon, Materials Discovery, **2016**, 6, 9-16

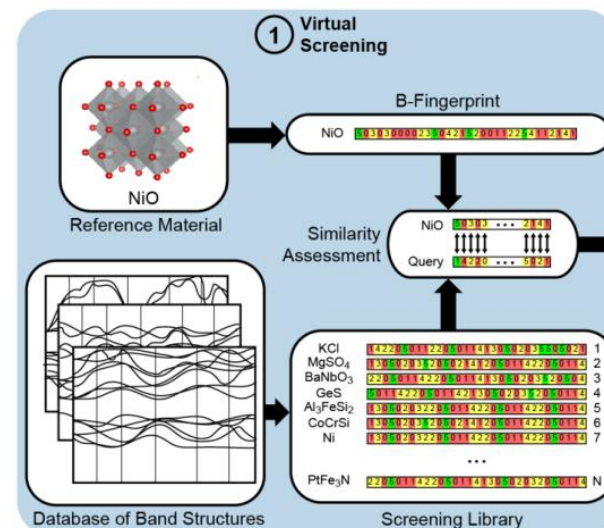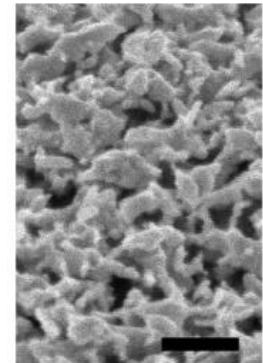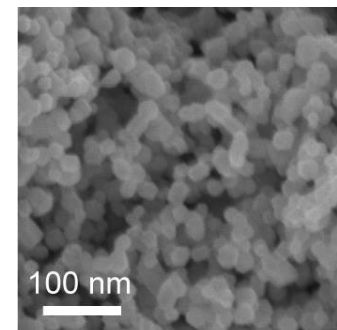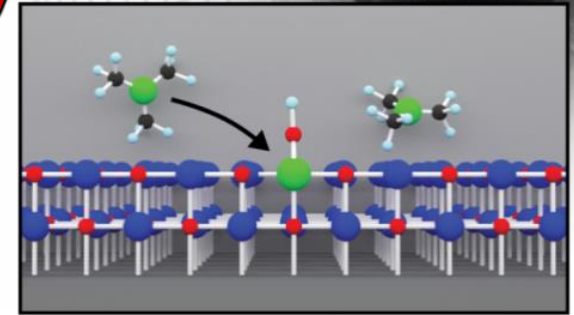# Summary of Materials Informatics: Supporting Experimental Discovery



PbTiO3 is identified as a new photocathode material.
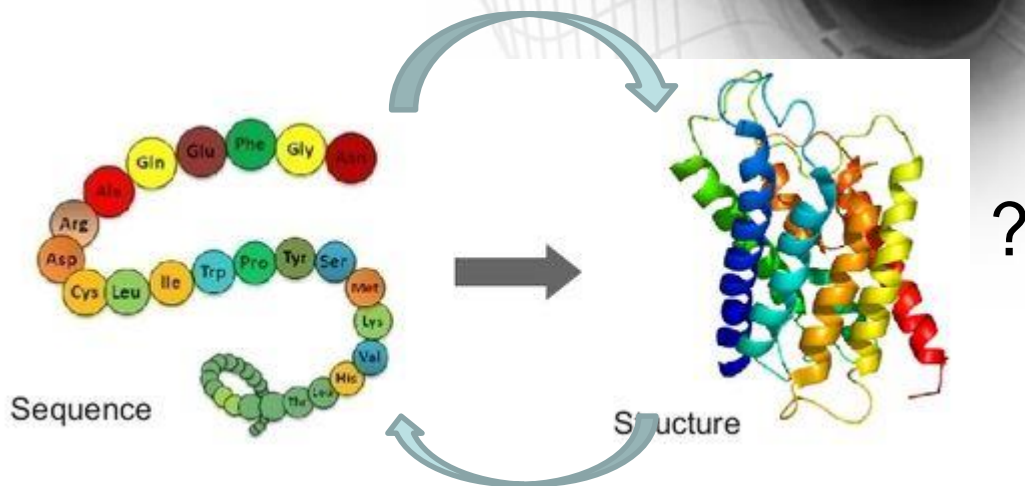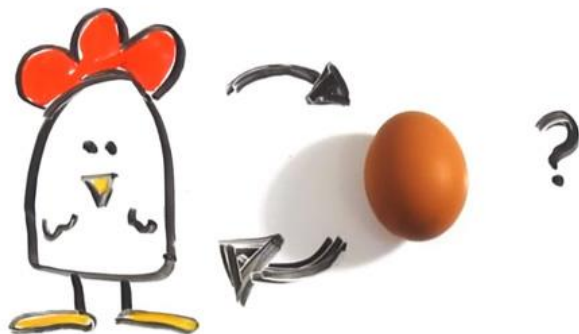
Successful experimental validation

Record fill factors of >50

First fully aqueous DSSC device

Currently, device performance is low; possible improvement by designing a new dye



100 nm



① Virtual Screening

B-Fingerprint

NiO

Reference Material

Similarity Assessment

NiO

Query

KCl
MgSO4
BaNbO3
GeS
Al3FeSi2
CoCrSi
Ni

PtFe3N

Database of Band Structures

Screening Library

# The eternal philosophical question:
# Which came first?



R or R$^2$    ?

In the beginning was the Word…
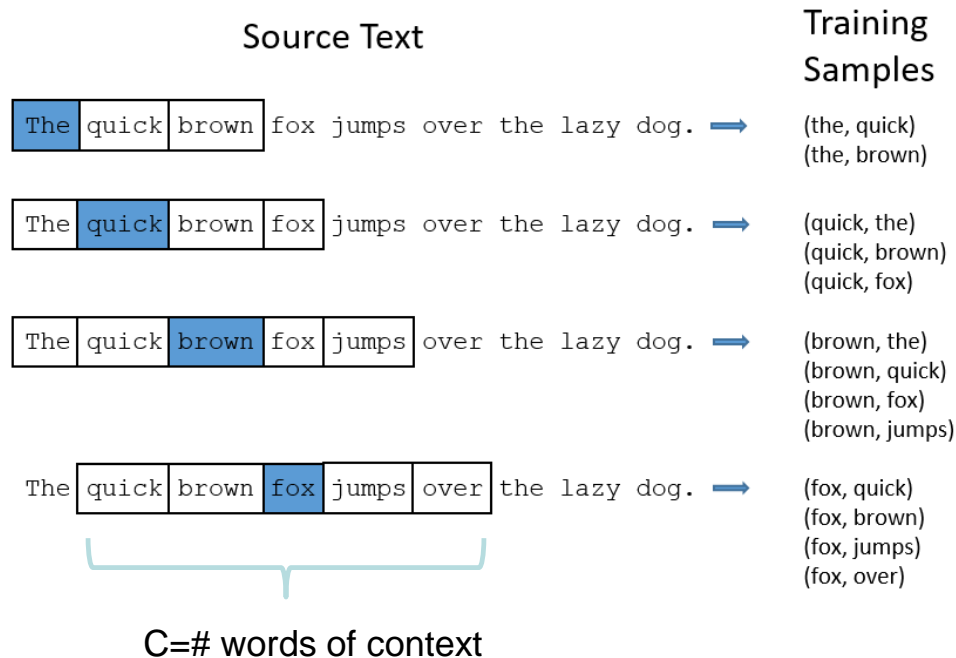
And the Word was…          **embedded**

(freely adopted from the Gospel of John)

"You should know a word by the company it keeps"
J.R.Firth 1957

British linguist; formulated the notion of the "context-dependent nature of meaning"

# Learning semantic context with Word2Vec



Source Text | Training Samples

The quick brown fox jumps over the lazy dog. ➡ (the, quick) (the, brown)

The quick brown fox jumps over the lazy dog. ➡ (quick, the) (quick, brown) (quick, fox)

The quick brown fox jumps over the lazy dog. ➡ (brown, the) (brown, quick) (brown, fox) (brown, jumps)

The quick brown fox jumps over the lazy dog. ➡ (fox, quick) (fox, brown) (fox, jumps) (fox, over)

C=# words of context
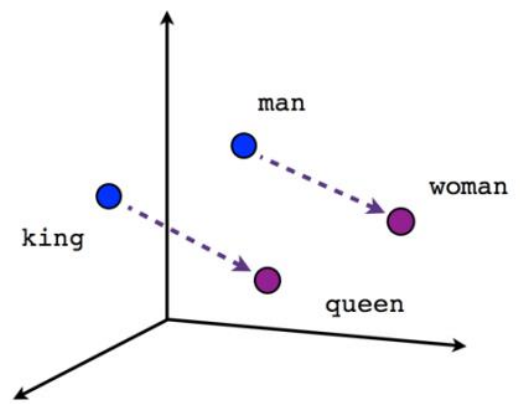
Can be used to learn:

CBOW*:*
- *Pr(word_k|words_context)*

Skip-Gram:
- *Pr(words_context|word_k)*

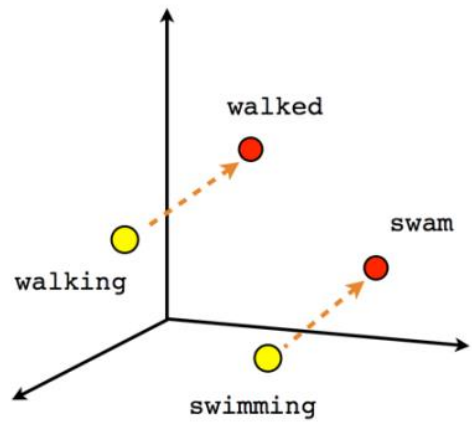- Mikolov, Tomas; et al. "Efficient Estimation of Word Representations in Vector Space". arXiv:1301.3781

Word2Vec Images courtesy of Chris McCormick:
http://mccormickml.com/2016/04/19/word2vec-tutorial-the-skip-gram-model/
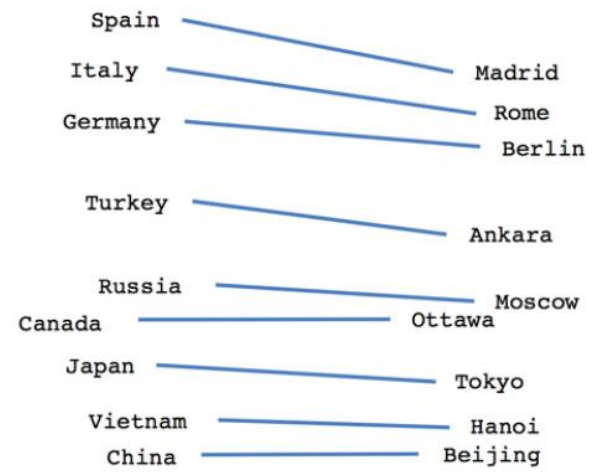
# Word embedding and similarity in the semantic space



Male-Female

Verb tense

Country-Capital

# SMILES are words that uniquely describe sentence-molecules!

Aspirin, also known as `O=C(C)Oc1ccccc1C(=O)O`, is a medication used to treat pain, fever, and inflammation.

`CC(=C)C1(C=CC(=O)O1)O` is a mycotoxin that is produced by Aspergillus flavus and *Penicillium roqueforti* mold.

| C O M P O U N D S | | | A C T I V I T Y |
|---|---|---|---|
| | O=C(C)Oc1ccccc1C(=O)O | Active | 1 |
| | CCOc1cc(C)ccc1OCC=CF | Inactive | 0 |
| | COc1ccccc1OCCO | Inactive | 0 |
| | CC(N)Sc1ccc(Cl)nc1 | Inactive | 0 |
| | COC(=O)NCc1ccccc1Cl | Active | 1 |

# ReLeaSE* design principles: learning and exploiting structural linguistics of SMILES notation

- SMILES notations reflect rules of Chemistry

- SMILES notation embeds linguistic rules

- Neural nets could learn both of the above types of rules

- This knowledge can be transformed into the generation of new SMILES corresponding to novel chemically feasible molecules (generative model)

- One can build QSAR models based solely on SMILES notation (predictive model)

- QSAR models can be used as a reward function for reinforcement learning to bias the design of novel libraries
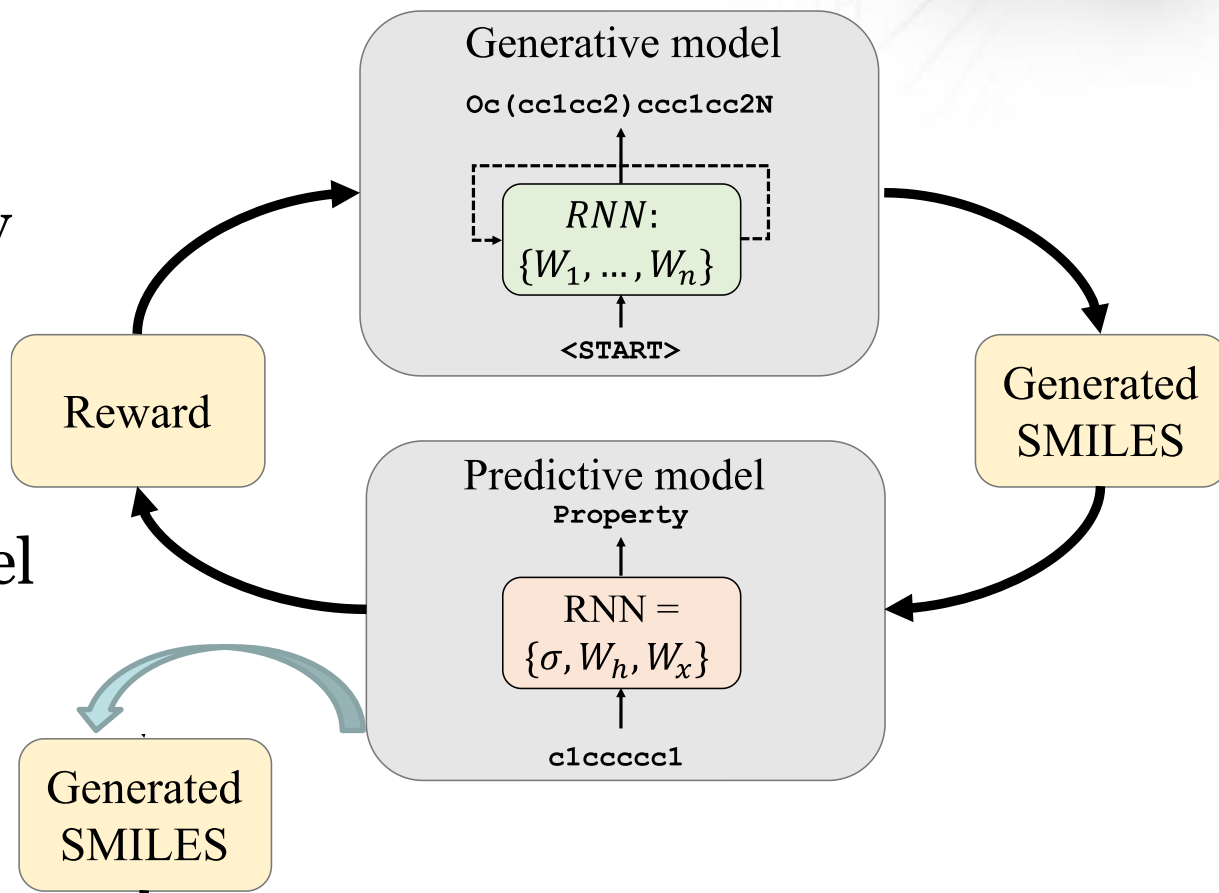
# Design of the ReLeaSE* method
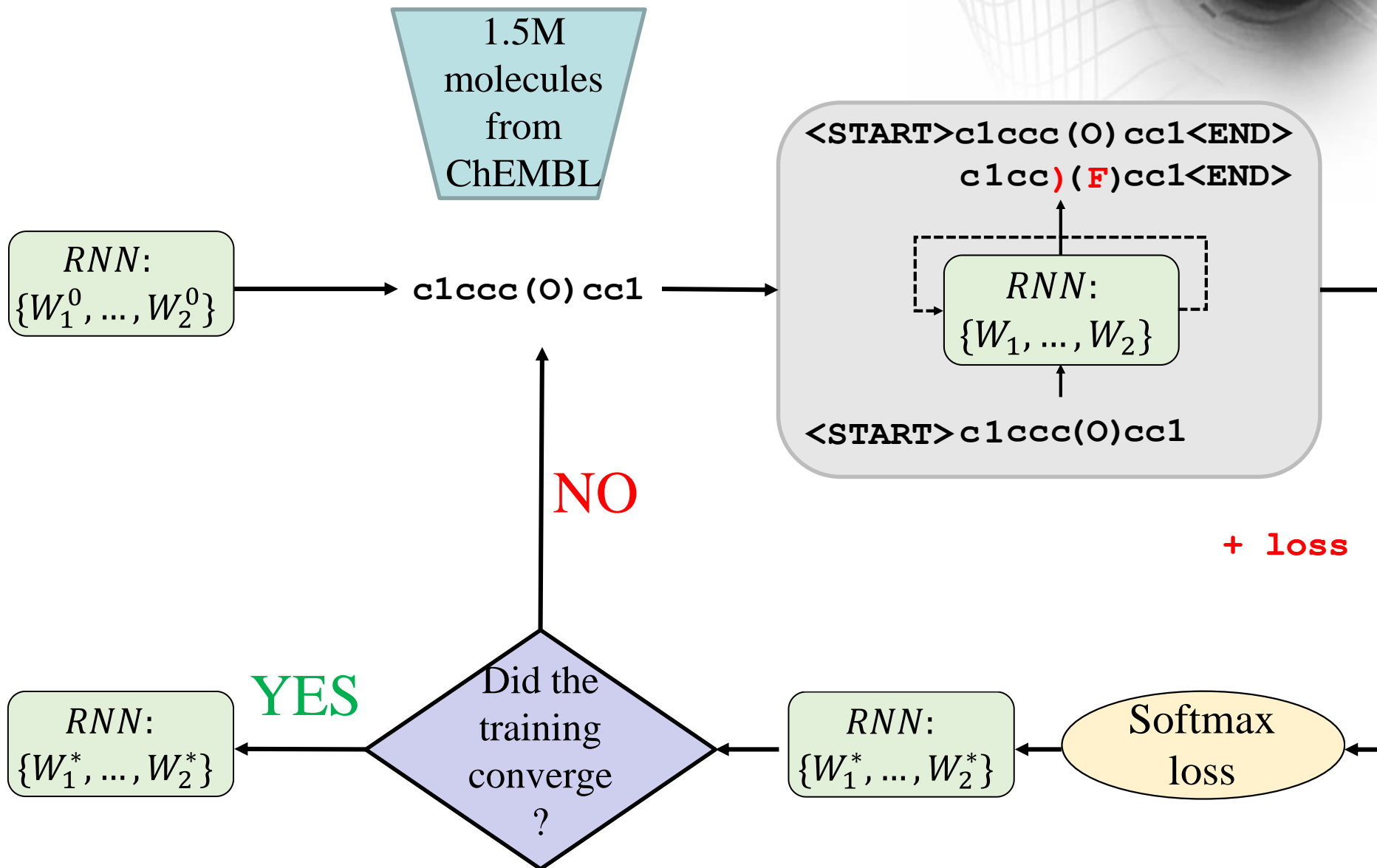## (Reinforcement Learning for Structural Evolution)

Elements of the thought cycle (molecules->models-molecules):

- Generate chemically feasible SMILES

- Develop SMILES-based QSAR model

- Employ QSAR model to bias library generation

- Produce new SMILES

**Generative model**

$Oc(cc1cc2)ccc1cc2N$

$RNN$:
$\{W_1, \ldots, W_n\}$

$<START>$

**Reward**

**Generated SMILES**

**Predictive model**

Property

$RNN = \{\sigma, W_h, W_x\}$
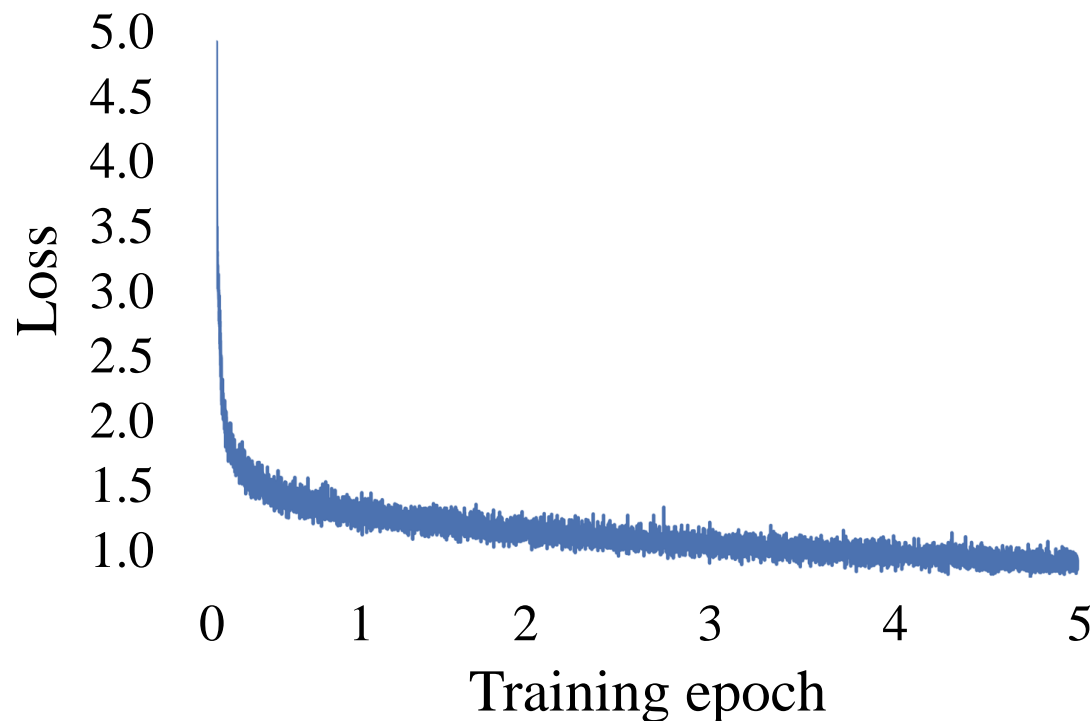
$c1ccccc1$

**Generated SMILES**

*Popova, Mariya, Olexandr Isayev, and Alexander Tropsha. "Deep reinforcement learning for de-novo drug design."
*arXiv preprint arXiv:1711.10907* (2017); Science Advances (in press).

# Generative model: training mode

# Generative model: training mode

- Training continues until convergence

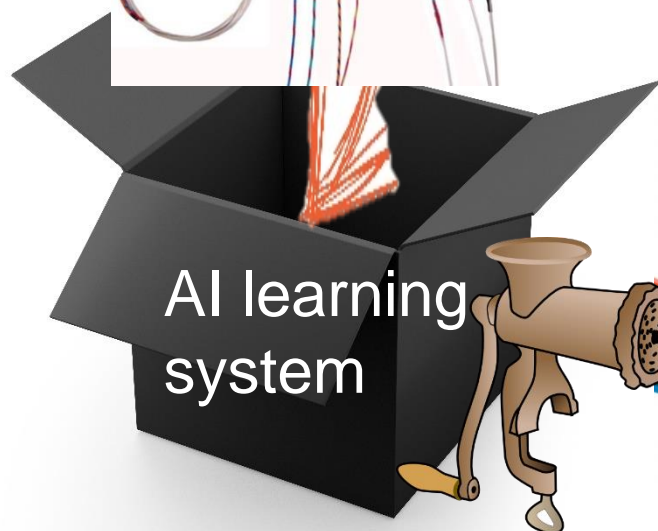- Every SMILES from ChEMBL is used as training example ~ 3-5 times
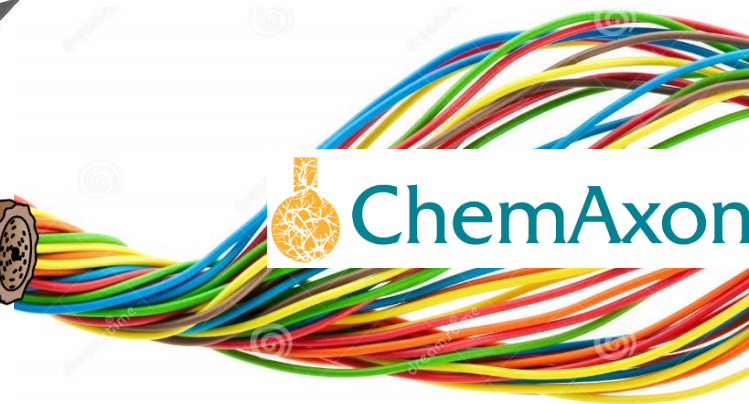
# Are we making legitimate Smiles?



PubChem / ChEMBL

Smiles strings

AI learning system

ChemAxon

95% Valid Chemically-feasible molecules

# Smile-ification of QSAR!

**COMPOUNDS**

O=C(C)Oc1ccccc1C(=O)O
CCOc1cc(C)ccc1OCC=CF
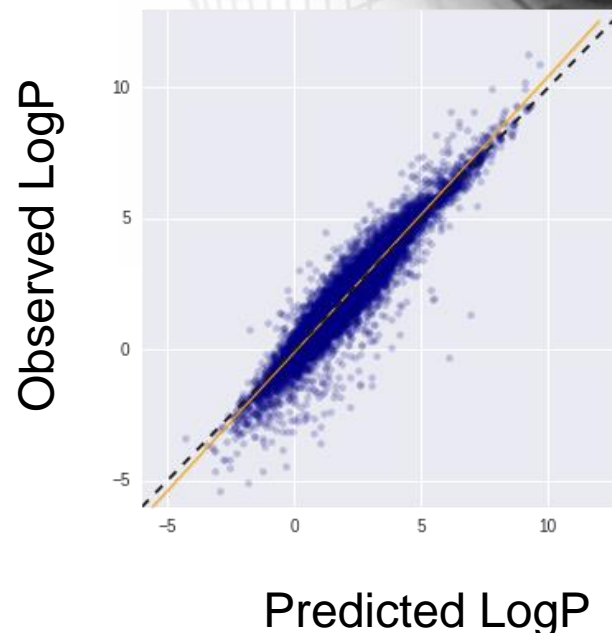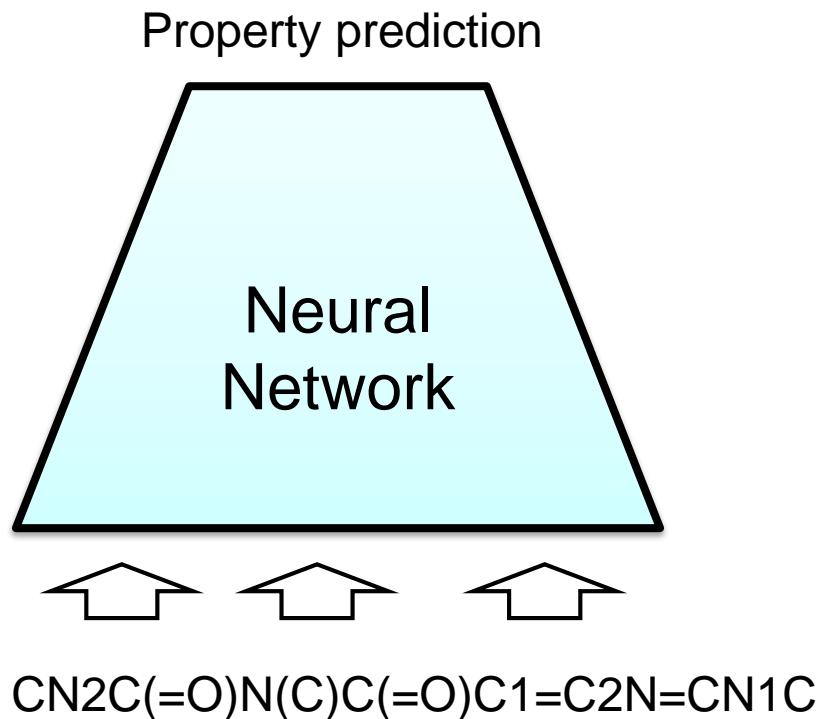COc1ccccc1OCCO
CC(N)Sc1ccc(Cl)nc1
COC(=O)NCc1ccccc1Cl

QSAR

**ACTIVITY**

0.531
1.299
0.946
-0.218
0.017

Quantitative <u>Smiles</u> – Activity Relationships

# QSAR modeling using Smiles strings only*

Property prediction

Neural Network

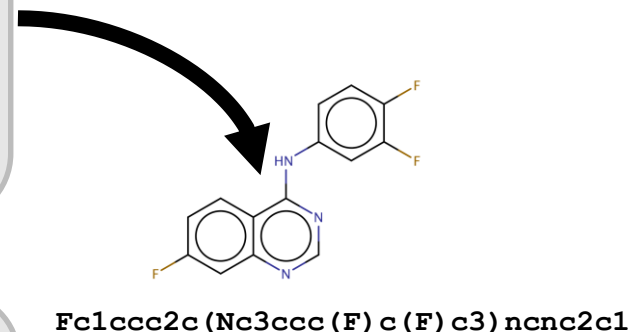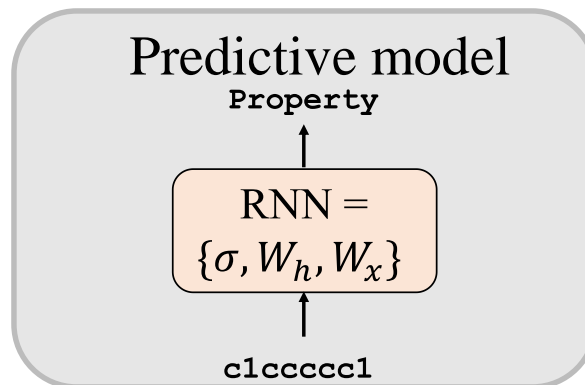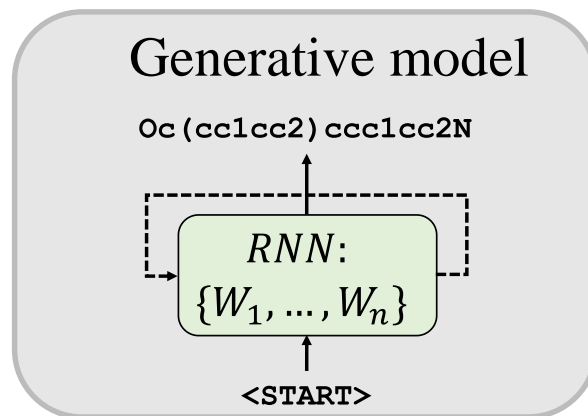CN2C(=O)N(C)C(=O)C1=C2N=CN1C



Observed LogP

Predicted LogP

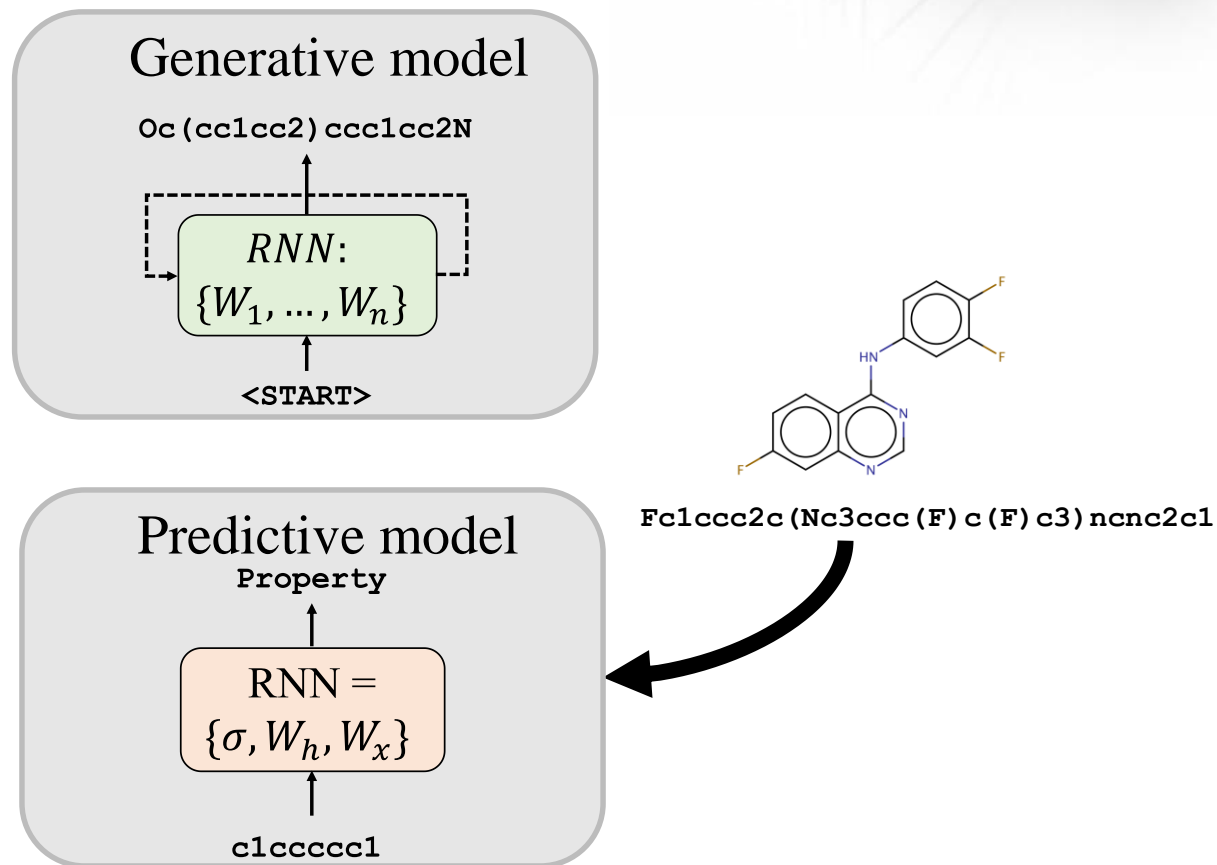|  | 5CV RF model with DRAGON7 Descriptors | 5CV NN model with SMILES directly |
|---|---|---|
| RMSE: | 0.57 | 0.53 |
| MAE: | 0.37 | 0.35 |
| $R^2_{ext}$: | 0.90 | 0.91 |

*LogP data for ~16K molecules from PHYSPROP (srcinc.com), Toxcast Dashboard (https://comptox.epa.gov/dashboard), and others.
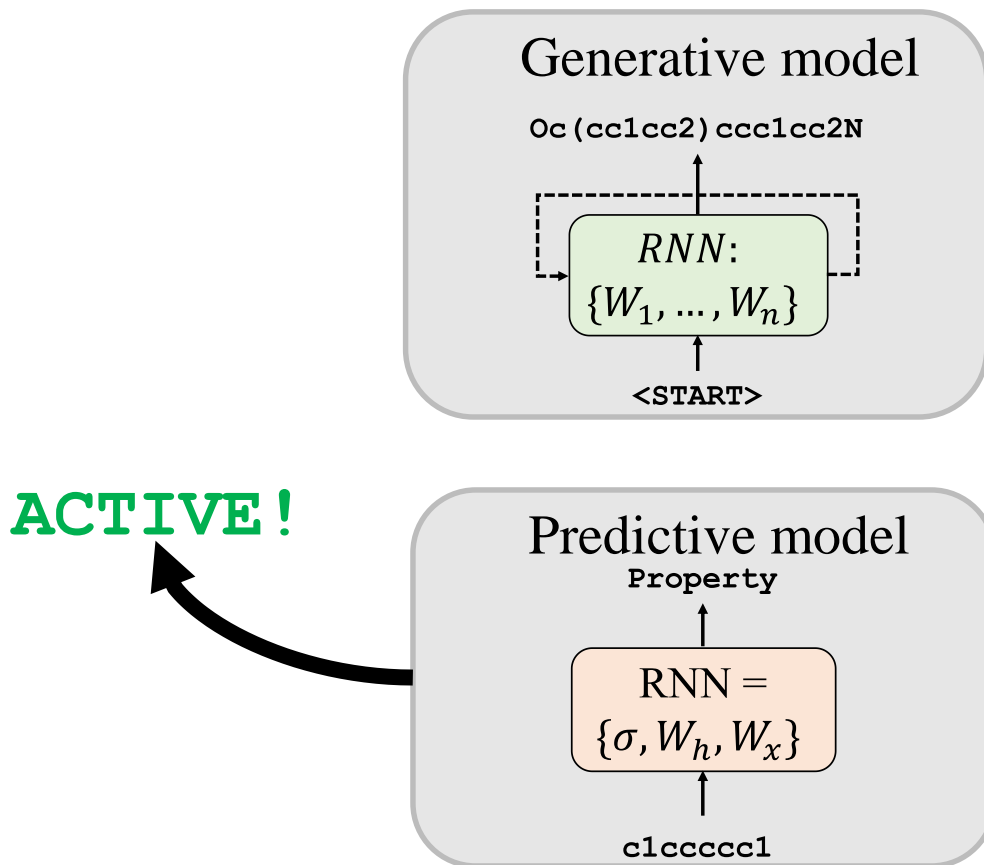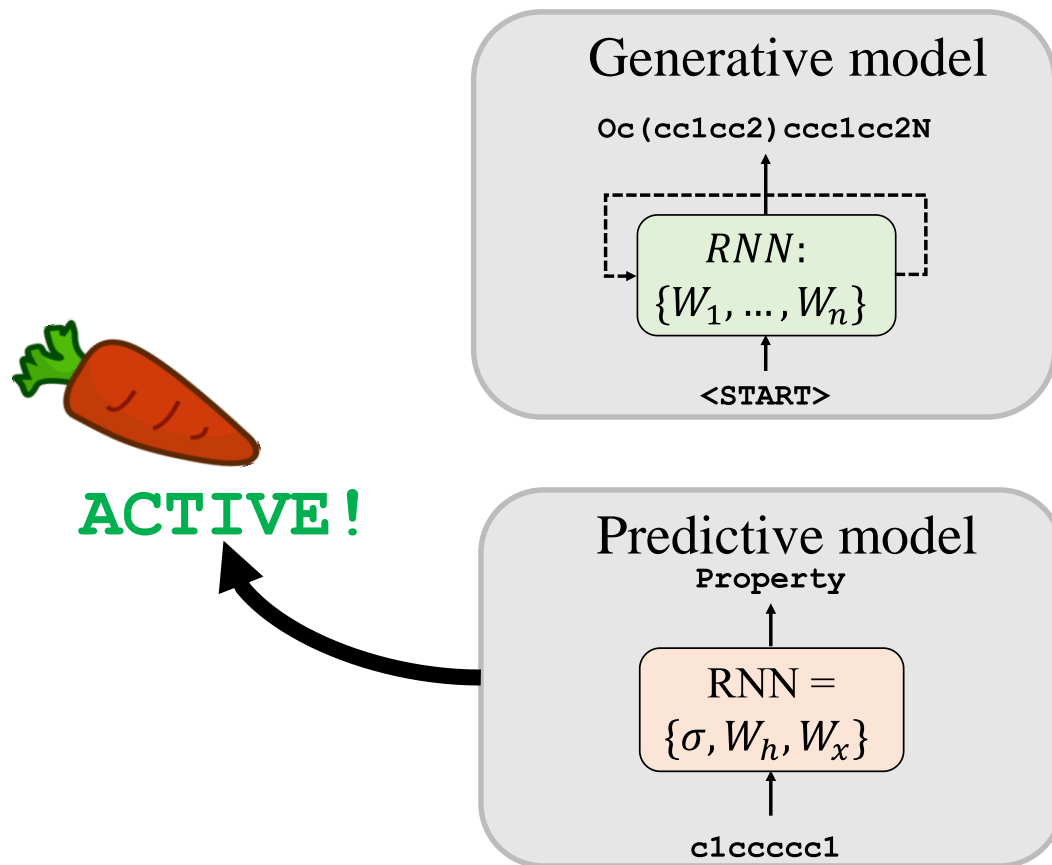
# Reinforcement learning for chemical design

# Reinforcement learning for chemical design

# Reinforcement learning for chemical design

# Reinforcement learning for chemical design

# Reinforcement learning for chemical design



Generative model

`Oc(cc1cc2)ccc1cc2N`

$RNN$:
$\{W_1, \dots, W_n\}$

`<START>`

**ACTIVE!**

Predictive model

`Property`

$RNN =$
$\{\sigma, W_h, W_x\}$

`c1ccccc1`

# Reinforcement learning for chemical design



Generative model

Predictive model

**Property**

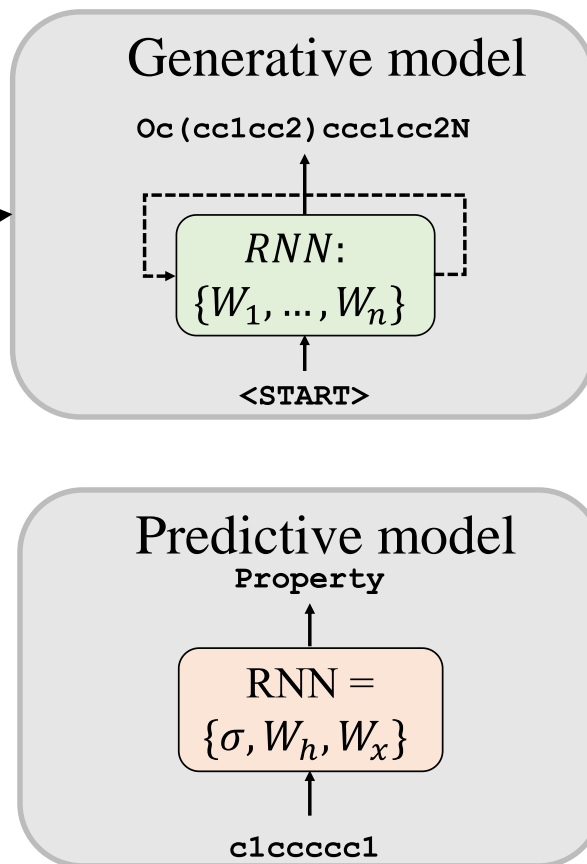$$RNN = \{\sigma, W_h, W_x\}$$

`c1ccccc1`

# Reinforcement learning for chemical design

# Reinforcement learning for chemical design

Generative model

`Oc(cc1cc2)ccc1cc2N`

$RNN$:
$\{W_1, \ldots, W_n\}$

`<START>`

Predictive model

`Property`

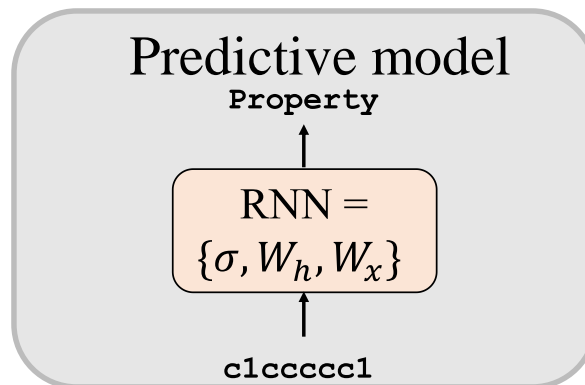$RNN = \{\sigma, W_h, W_x\}$

`c1ccccc1`
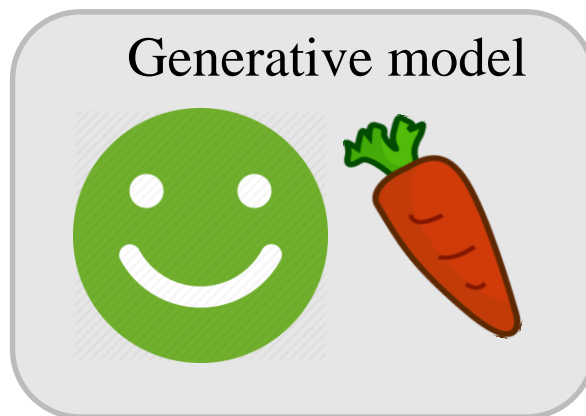
`FC(F)COc1ccc2c(Nc3ccc(Cl)c(Cl)c3)ncnc2c1`

# Reinforcement learning for chemical design

# Reinforcement learning for chemical design

# Reinforcement learning for chemical design

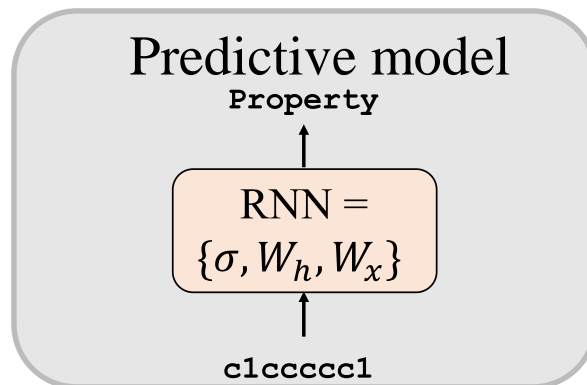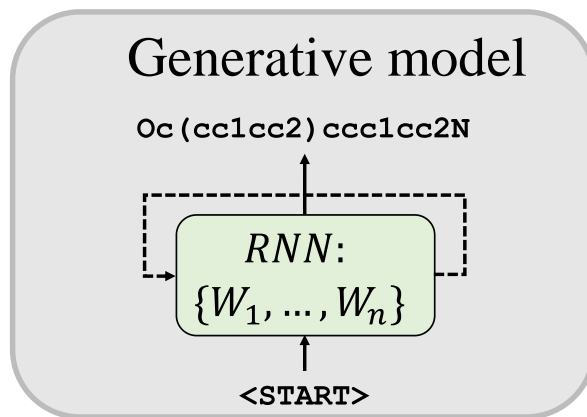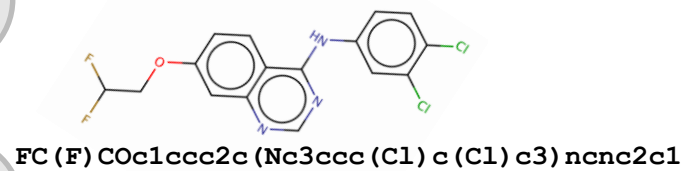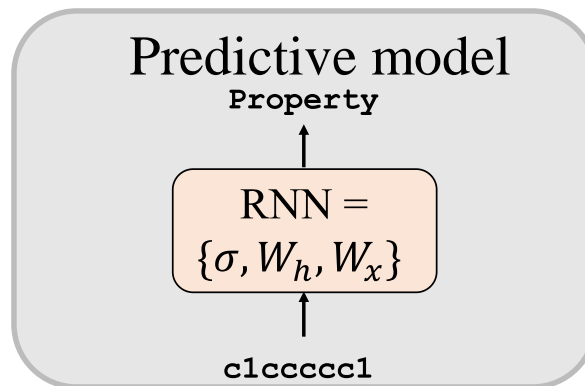Generative model

`Oc(cc1cc2)ccc1cc2N`

$RNN$:
$\{W_1, ..., W_n\}$

`<START>`

**INACTIVE!**

Predictive model

`Property`

$RNN =$
$\{\sigma, W_h, W_x\}$

`c1ccccc1`

# Reinforcement learning for chemical design



Generative model

Predictive model

Property

$$RNN = \{\sigma, W_h, W_x\}$$

c1ccccc1

# Technical details

- Models were trained on Nvidia Titan X and Titan V GPUs

- Training the generative model on ChEMBL took ~ 25 days

- Training of predictive models took ~ 2 hours

- Biasing the generative model with reinforcement learning for one property ~ 1 day

- Generative model produces 1000 compounds per minute

# Results: Synthetic accessibility score* of the designed libraries



Synthetic accessibility score

*Ertl, Peter, and Ansgar Schuffenhauer. "Estimation of synthetic accessibility score of drug-like molecules based on molecular complexity and fragment contributions." *Journal of cheminformatics* 1.1 (2009): 8.

# PoC: Structural Bias

A: increase in number of substituents

Reward increase

B: increase in number of benzene rings

# Results: Biasing target properties in the designed libraries



Melting temperature ($T_m$), ºC

Number of substituents

# Results: Biasing target properties in the designed libraries

# Target predictions for generated compounds using SEA*



| Query | Target Key | Target Name | Description | P-Value | MaxTC |
|---|---|---|---|---|---|
| | NPM_HUMAN+5 | NPM1 | Nucleophosmin | 3.118e-74 | 0.49 |
| | CCNH_HUMAN+5 | CCNH | Cyclin-H | 2.571e-32 | 0.38 |
| | PAK1_HUMAN+5 | PAK1 | Serine/threonine-protein kinase PAK 1 | 5.277e-24 | 0.39 |
| | ALK_HUMAN+5 | ALK | ALK tyrosine kinase receptor | 3.714e-23 | 0.54 |
| | JAK2_HUMAN+5 | JAK2 | Tyrosine-protein kinase JAK2 | 1.136e-21 | 0.61 |
| | INSR_HUMAN+5 | INSR | Insulin receptor | 2.36e-17 | 0.54 |
| | CCNB1_HUMAN+5 | CCNB1 | G2/mitotic-specific cyclin-B1 | 2.22e-16 | 0.38 |

*Keiser MJ, Roth BL, Armbruster BN, Ernsberger P, Irwin JJ, Shoichet BK. Relating protein pharmacology by ligand chemistry. *Nat Biotech* **25** (2), 197-206 (2007).

# Target predictions for generated compounds using SEA*



| Query | Target Key | Target Name | Description | P-Value | MaxTC |
|---|---|---|---|---|---|
| | EGFR_HUMAN+5 | EGFR | Epidermal growth factor receptor | 8.688e-244 | 0.61 |
| | ERBB2_HUMAN+5 | ERBB2 | Receptor tyrosine-protein kinase erbB-2 | 8.544e-169 | 0.55 |
| | ERBB2_RAT+5 | Erbb2 | Receptor tyrosine-protein kinase erbB-2 | 5.893e-87 | 0.42 |
| | VGFR2_HUMAN+5 | KDR | Vascular endothelial growth factor receptor 2 | 6.294e-65 | 0.58 |
| | ERBB4_HUMAN+5 | ERBB4 | Receptor tyrosine-protein kinase erbB-4 | 1.354e-64 | 0.49 |

*Keiser MJ, Roth BL, Armbruster BN, Ernsberger P, Irwin JJ, Shoichet BK. Relating protein pharmacology by ligand chemistry. *Nat Biotech* **25** (2), 197-206 (2007).

# Results: analysis of similarity

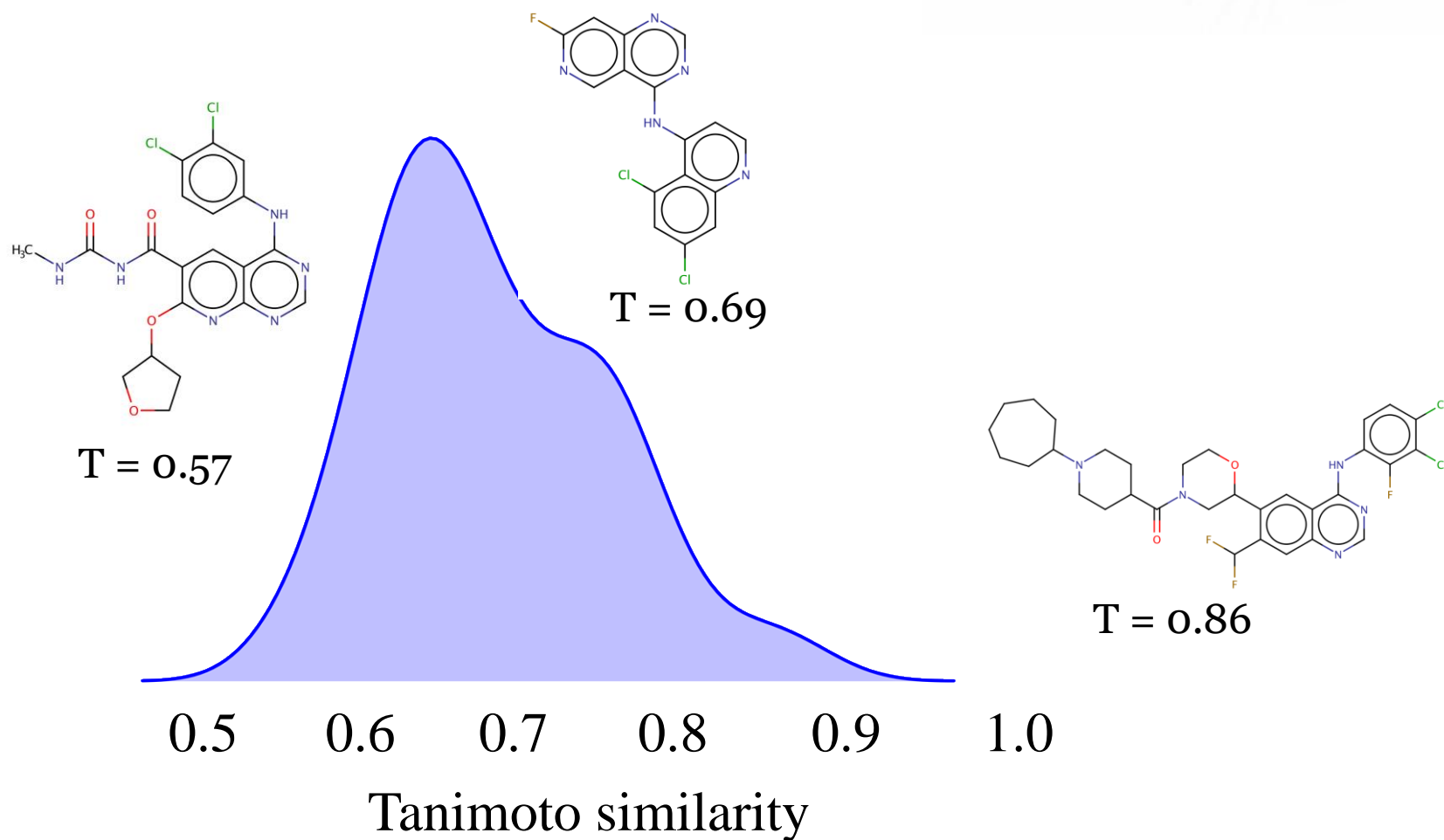Distribution of Tanimoto similarity to the nearest neighbor in training dataset for compounds predicted to be active for EGFR by consensus of QSAR models:



T = 0.69

T = 0.57

T = 0.86

0.5    0.6    0.7    0.8    0.9    1.0

Tanimoto similarity

# Model visualization for putative JAK2 inhibitors (projection using t-SNE)



ZINC469992
pIC50 = 8.23

ZINC19982368
pIC50 = 8.64

ZINC2876515
pIC50 = 8.39

ZINC66347860
pIC50 = 3.31

pIC50 = 0.63

ZINC3549031
pIC50 = 3.76

pIC50 = 10.37

# Summary

- AI methods coupled with SMILES representation (only!) afford biased library generation

- The system naturally embeds reinforcement learning to produce novel structures with the desired property

- The system can be tuned to bias libraries towards specific property ranges

- Next phase is experimental validation of hits

# Summary of recent AI-based studies on chemical library design

| Molecular representations | Generative models | Method of biasing generated compounds |
|---|---|---|
| • Fingerprints<br>• SMILES<br>• Graphs | • Autoencoders<br>• Generative adversarial models<br>• Recurrent neural networks<br>• Convolutional neural networks | • None<br>• Latent space optimization<br>• Fine-tuning on small subset of molecules with the desired property<br>• Reinforcement Learning |

# An example of experimental validation of AI-based models*

- First training on large dataset
- Then fine-tuning on small subset of active compounds
- "These observations corroborate the ability of the generative AI model to **produce novel chemical entities within the training data domain**".
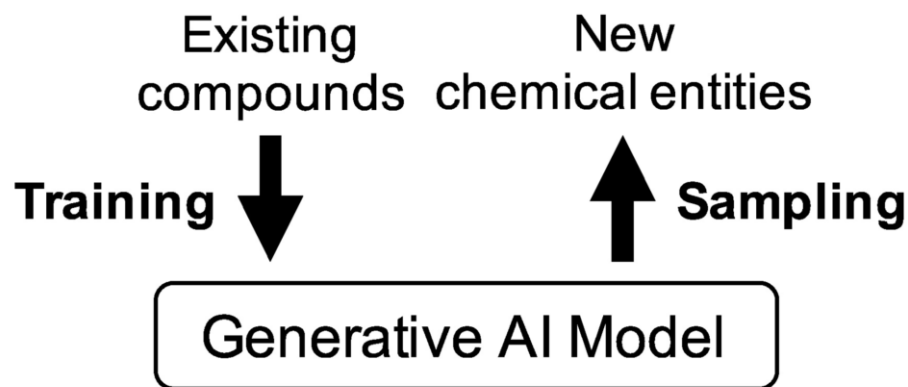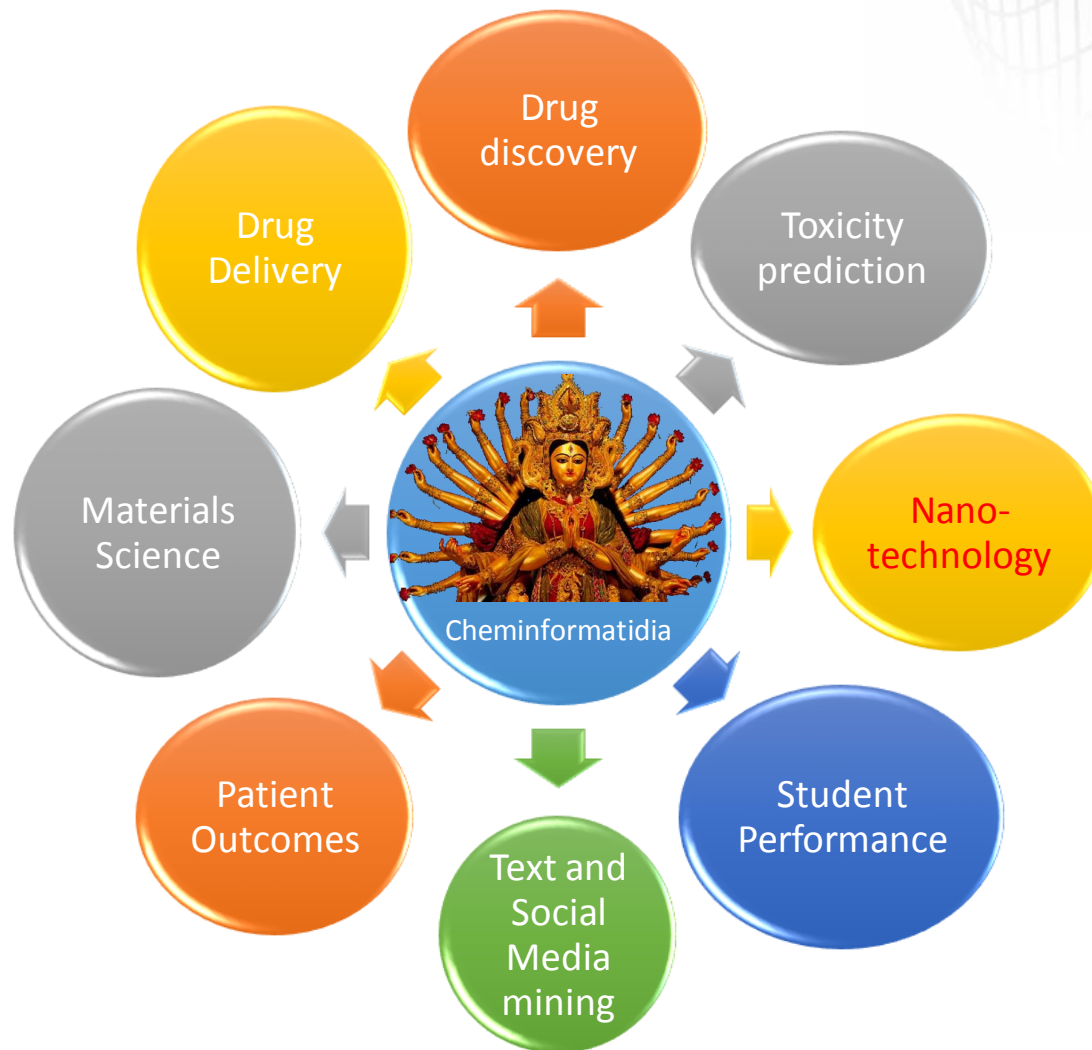
**Table 1.** *In vitro* activity of designs **1–5** on RXRs and PPARs ($EC_{50}$ values $\pm$ SEM [µM]; $n = 2$ (when inactive) or 4 (when active) independent experiments in duplicates; *inactive*, no statistically significant reporter transactivation at a compound concentration of 30 µM).

| Compound no. | RXRα | RXRβ | RXRγ | PPARα | PPARγ | PPARδ |
|---|---|---|---|---|---|---|
| 1 | $0.13 \pm 0.01$ | $1.1 \pm 0.3$ | $0.06 \pm 0.02$ | *inactive* | $2.3 \pm 0.2$ | *inactive* |
| 2 | $13.0 \pm 0.1$ | $9 \pm 2$ | $8.0 \pm 0.7$ | *inactive* | $2.8 \pm 0.3$ | *inactive* |
| 3 | *inactive* | *inactive* | *inactive* | $4.0 \pm 1.0$ | $10.1 \pm 0.3$ | *inactive* |
| 4 | *inactive* | *inactive* | *inactive* | *inactive* | $9 \pm 3$ | $14 \pm 2$ |
| 5 | *inactive* | *inactive* | *inactive* | *inactive* | *inactive* | *inactive* |
| reference agonists[a] | $0.033 \pm 0.002$ | $0.024 \pm 0.004$ | $0.025 \pm 0.002$ | $0.006 \pm 0.002$ | $0.6 \pm 0.1$ | $0.5 \pm 0.1$ |

[a] Reference agonists, literature data: bexarotene[17] for RXRs, GW7647[18] for PPARα, pioglitazone[19] for PPARγ, L165,041[19] for PPARδ

* D. Merk, L. Friedrich, F. Grisoni, G. Schneider, *Mol. Inf.* **2018**, *37*, 1700153.

# Many virtues of Cheminformatics

# Acknowledgements

**Principal Investigator**
**Alexander Tropsha**

**Research Professors**
**Alexander Golbraikh**
**Olexander Isayev**
**Eugene Muratov**

**Postdoctoral Fellows**
**Vinicius Alves**
**Stephen Capuzzi**
**Joyce Borba**

**Graduate students**
**Sherif Faraq**
**Kyle Bowers**
**Maria Popova**
**Andrew Thieme**

- **Duke University**
  – Stefano Curtarolo
  – Corey Oses
- **UNC Chemistry**
  – Jim Cahoon
  – Taylor Moot
  – Aaron Taggart

67