# CDD WEBINAR

## Generative Chemistry, Deep Learning and Traditional Models
### Moderated by Dr. Eric Gifford

**LIVE**

**Thursday, September 26th 2024**
**8:00 AM (PT) | 11:00 AM (ET) | 16:00 (BST)**

**Eric Gifford, PhD.**
Business Development Consultant
Collaborative Drug Discovery

**Pat Walters, PhD.**
Chief Data Officer
Relay Therapeutics

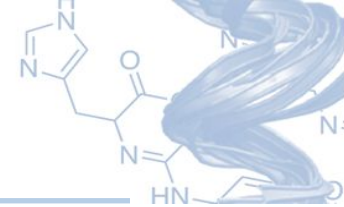**Greg Landrum, PhD.**
Senior Scientist
ETH Zürich

**Peter Gedeck, PhD.**
Senior Informatician
Collaborative Drug Discovery

# Do you have a question to ask our panel?

Open the ZOOM Q&A
type in your question during the webinar



*We will reserve time and answer as many questions as we can at the end*

CDD WEBINAR

# Generative Chemistry, Deep Learning and Traditional Models
## Moderated by Dr. Eric Gifford

**LIVE**

**Thursday, September 26th 2024**
**8:00 AM (PT) | 11:00 AM (ET) | 16:00 (BST)**

Eric Gifford, PhD.
Business Development
Consultant
Collaborative Drug Discovery

Pat Walters, PhD.
Chief Data Officer
Relay Therapeutics

Greg Landrum, PhD.
Senior Scientist
ETH Zürich

Peter Gedeck, PhD.
Senior Informatician
Collaborative Drug
Discovery

# What does our field need to move forward?

## *(Some topics for today's discussion)*

◆ **Comparable Data** ◆

◆ **Open-source tool kits and commercial software** ◆

◆ **Robust Models & AI** ◆

# Pop quiz

Suppose you're doing machine learning and you need some bioactivity (or ADME) data sets to work with.

You want/need to use public data, so you start with ChEMBL

You notice that you can increase the size of the data sets available to you *significantly* if you combine results taken from different assays but run on the same target

Is it safe to do this?

# Unfortunately common in our field

"I really need more data, so I'm just going to combine the assays. How bad can it possibly be?"

# Aside: the Cheng-Prusoff equation

Relates IC50 to Ki for competitive inhibition assays:

$Ki = IC50/(1+ [S]/Km)$

where:

[S] = concentration of the substrate

Km = the affinity constant of the substrate

This does *not* magically let you convert IC50 (generally unsafe to compare across assays) to Ki (safe to compare across assays): you need to know both [S] and Km in order to use it, and if you know [S] and Km then you *already* know whether or not it's safe to compare

A fun take on the relationship between Ki and IC50:

https://krhornberger.substack.com/p/tweetorial-ic50-vs-ki

# It's pretty bad.



IC50, only activity curation

1599 assay pairs. 50385 points.
R2=0.32, Spearman R=0.68 MAE=0.51

Fraction > 0.3: 0.65
Fraction > 1.0: 0.28

# Having a Realistic Dynamic Range is Important



**Biogen 3.5 logs**
**Delaney 13.2 logs**

dataset
— Delaney
— Biogen

**The Delaney solubility dataset is a component of MoleculeNet**

# Focus on Simple, Consistent, Well-Defined Endpoints



**Physiology**

- BBBP
- Tox21
- ToxCast
- SIDER
- ClinTox

**SIDER Categories**

Hepatobiliary disorders
Metabolism and nutrition disorders
Product issues
Eye disorders
Investigations
Musculoskeletal and connective tissue disorders
Gastrointestinal disorders
Social circumstances
Immune system disorders
Reproductive system and breast disorders
Neoplasms benign, malignant and unspecified (incl cysts and polyps)
General disorders and administration site conditions
Endocrine disorders
Surgical and medical procedures
Vascular disorders
Blood and lymphatic system disorders
Skin and subcutaneous tissue disorders
Congenital, familial and genetic disorders
Infections and infestations
Respiratory, thoracic and mediastinal disorders
Psychiatric disorders
Renal and urinary disorders
Pregnancy, puerperium and perinatal conditions
Ear and labyrinth disorders
Cardiac disorders
Nervous system disorders
Injury, poisoning and procedural complications

**CDD Chemistry Engine**
- Backend uses RDKit
- Frontend uses RDKit and WebMolKit (Alex Clark)

**CDD contributions to RDKit:**
- Representation of atropisomers

**Ketcher (open source structure editor):**
- Work with *epam* to implement our requirements
- Focus on adding macromolecule functionality (non-natural peptides, ADC, …)

**CDD Visualization**
- Not open source, but scientists can use CDD Visualization for free



https://www.collaborativedrug.com/scientific-data-visualization-software

Greg Landrum
@dr_greg_landrum@sciencemastodon.com
@greg_landrum.bsky.social

## Usage in Commercial Tools

- Amazon Web Services
- Collaborative Drug Discovery
- Cresset Software
- Dalke Scientific Software
- Datagrok
- Glysade
- MedChemica
- NextMove Software
- Schrödinger
- SCM
- Wolfram Research

Disclaimer: this info is from public statements made by people from those companies.  I almost certainly have forgotten someone

## Adoption Measures

- Mailing lists: ~250 messages to rdkit-discuss from 2022.09 - 2023.08
- Google scholar: >2300 hits for "rdkit" in 2022, >2000 so far in 2023
- Searching github for "from rdkit import Chem" returns >27000 code results
- Each of the last nine in-person UGMs at capacity with 40-150 attendees

# Models, Tools, Technologies for Comparisons

# How can we compare different modeling methods?

# Don't Compare Mean Performance Across Cross-Validation Folds



| Method | Mean $R^2$ |
|:------:|:----------:|
| 1 | 0.46 |
| 2 | **0.48** |

# We Have Distributions Across Folds



p=0.06
d=0.55
medium effect

There are well established statistical tests for comparing distributions

**Thoughts on when it is best to use different methods:
AI, Deep Learning, Generative QSAR and/or Structure-Based
regression vs classification models?**

**What are the key differentiators?**

**CDD Deep Learning model**
- Create *unique* numerical representation of chemical structures (latent vector)
- Generate structures for a given latent vector

**Validation study showed**
- Latent vector encodes chemical structure
- Latent vector encodes structural relationships

Structure

↓

GCV

↓

Latent vector

Unique numerical representation

→ Generative Model

↓

Structure

Multiple representations possible

# Deep Learning: ChEMBL, SureChEMBL, New IP

**1) Generate Bioisosteres
2) Ultra-Fast Deep Learning Similarity searching**

**2. Similar compounds**
- **Sort by scaffold**
- **Calculate properties**

**3. Export compounds to a spreadsheet**

**4. New hits can be used for further analysis or purchasing**

**Complexity Simplified**
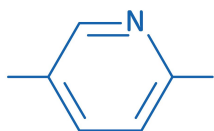
# Single Source of Truth

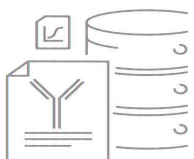## CDD VAULT®
### COLLABORATIVE DRUG DISCOVERY

**Inventory**
Keep track of samples, biologicals and compounds

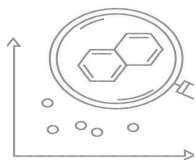**Activity**
Manage and analyze experimental data
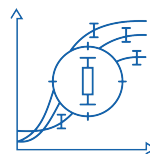
**Registration**
Store and organize your research data
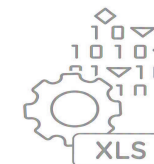
**Visualization**
Plot datasets and mine them
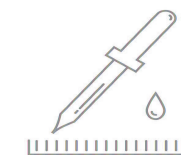
**Curves**
Generate, QC, and analyze results

**AI**
Computer aided design

**ELN**
Document all your research.

**Automation**
Connect data with our robust API, Parser, and Mapping Tools

**Assays**
Comparison of assays using standardized protocols

www.collaborativedrug.com          info@collaborativedrug.com

© All Rights Reserved Collaborative Drug Discovery