XXXI Symposium on Bioinformatics and Computer-Aided Drug Discovery (BCADD-2025)



On the use of Machine Learning for New Approach Methodologies

Tudor Oprea
Expert Systems Inc

toprea_at_expertsystems.inc https://bit.ly/oprea_ti

What is a Drug?



- A scientist's definition differs from a pharmacist's (perception matters)
- 2. Two major entities are interlinked in <u>DrugCentral</u>: the APIs, referred to as 'drugs' by chemists, biologists and other basic scientists, and pharmaceutical products, referred to as 'drugs' by patients, pharmacists, nurses and physicians, as well as other clinician scientists. (<u>source</u>)
- 3. The term 'drug' used here refers to therapeutic ingredients only [...] and does not include imaging agents, nutritional supplements, sunscreens or vaccines. (link)
- 4. ... but can you tell a drug when you see one?!
- 5. This decades-old question has developed into many papers that use ML to quantify "druglikeness", i.e., how similar are new agents to what others have made before (<u>link</u>)
- 6. Drug is not an intrinsic property of matter. Man giveth, Man taketh away.

Al vs Machine Learning



- "Machine learning is the study of computer algorithms that improve automatically through experience." -- Tom M. Mitchell
- 2. "Artificial intelligence is the science and engineering of making computers behave in ways that, until recently, we thought required human intelligence." -- Andrew Moore
- 3. First "AI"-term indexed in MeSH is from 1951 (KH Fletcher described a tortoise robot).

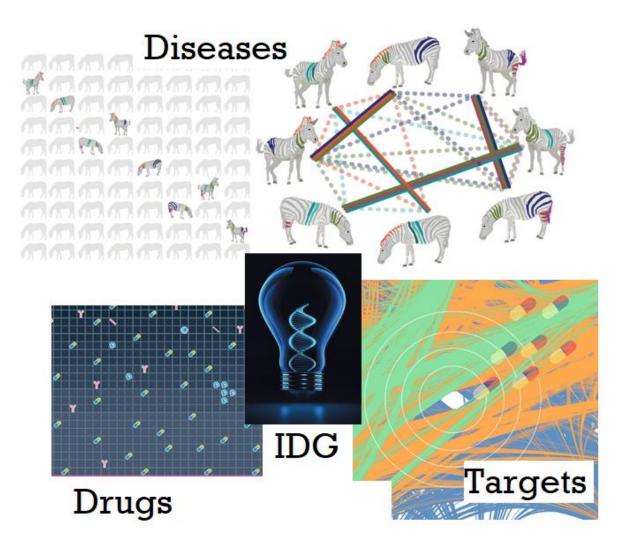
 Today we have speech recognition, self-driving cars, chatbots, image processing "AI", automated language translation, "generative AI" (movies, music, screenplays, etc).
- 4. Al systems have passed medical / legal board exams, are chess/GO champions, etc.
- 5. Unlike machine learning, AI is a *moving target*, as the definition of AI changes given rapid technological advancements. -- Roberto Iriondo
- 6. Except for the LLM section, everything in this talk describes ML models



Earlier Work:

Drug Targets

Three Pillars in Drug Discovery and Repurposing where AIML Contributes



Informatics, Data Science and Machine Learning ("AI") can be used as follows:

Diseases: EMR data extraction, nosology, ontology, & EMR-based AIML

Targets: drug target repurposing, selection & validation, phenotype associations AI/ML

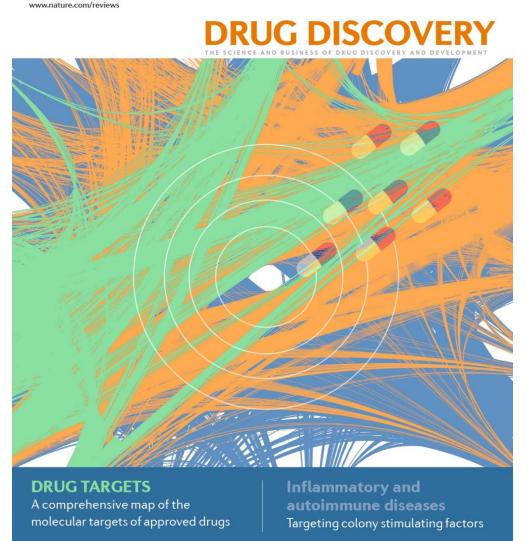
Drugs: Novel therapeutic modalities and drug repurposing using *in silico* methods

Drug discovery remains more art than science. Al/ML can help

Diseases image credit: Julie McMurry, Melissa Haendel (OHSU).
All other images credit: Nature Reviews Drug Discovery cover pages



A Comprehensive Map of Molecular Drug Targets



We curated 667 human genome-derived proteins and 226 pathogen-derived biomolecules through which 1,578 US FDA-approved drugs act.

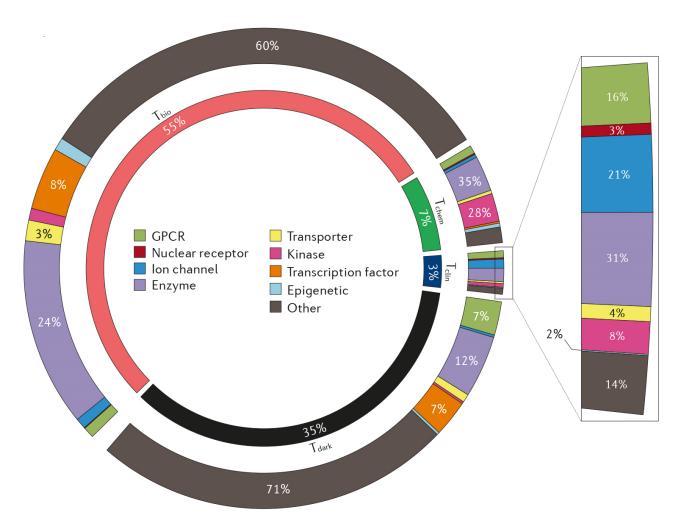
This set included 1004 oral drugs as well as 530 injectable drugs (approved through June 2016).

Data captured in DrugCentral (link)

* 2025 version submitted to NRDD



Knowledge-based Classification of Human Proteins



Most protein classification schemes are based on structural and functional criteria.

For therapeutic development, it is useful to understand how much and what types of data are available for a given protein, thereby highlighting well-studied and understudied targets.

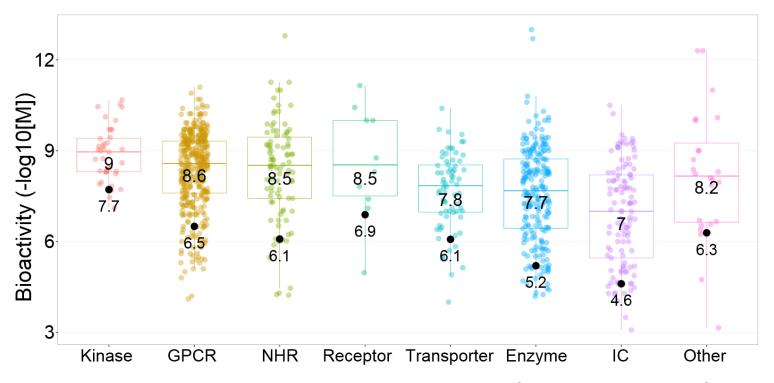
Tclin: Proteins annotated as drug targets **Tchem:** Proteins for which *potent* small molecules are known

Tbio: Proteins for which biology is better understood

Tdark: These proteins lack antibodies, publications or Gene RIFs



D-T Development Level 1



Bioactivities of approved drugs (by Target class)

ChEMBL: database of bioactive chemicals

https://www.ebi.ac.uk/chembl/

DrugCentral: online drug compendium

http://drugcentral.org/

Tclin proteins are associated with drug Mechanism of Action (MoA) – NRDD 2017

Tchem proteins have bioactivitis in ChEMBL and DrugCentral, + human curation for some targets

- Kinases: <= 30nM
- GPCRs: <= 100nM
- Nuclear Receptors: <= 100nM
- Ion Channels: <= 10μM
- Non-IDG Family Targets: <= 1μM



D-T Development Level 2

Tbio proteins lack small molecule annotation cf. Tchem criteria, and satisfy one of these criteria:

- protein is above the cutoff criteria for **Tdark**
- protein is annotated with a GO Molecular Function or Biological Process leaf term(s) with an Experimental Evidence code
- protein has confirmed <u>OMIM</u> phenotype(s)

Fractional paper count

PubMed score = $\sum_{j \in D} \frac{n_{ij}}{n_{.j}}$

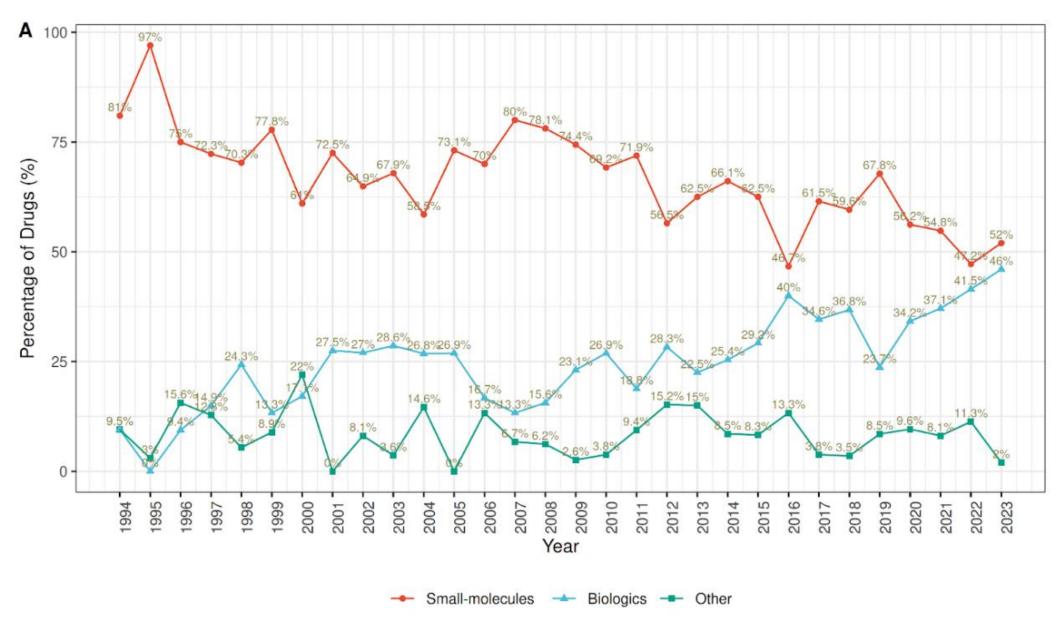
Tdark ("ignorome") have little information available, and satisfy these criteria:

- PubMed text-mining score from <u>Jensen Lab</u> < 5
- <= 3 Gene RIFs
- <= 50 Antibodies available according to <u>antibodypedia.com</u>



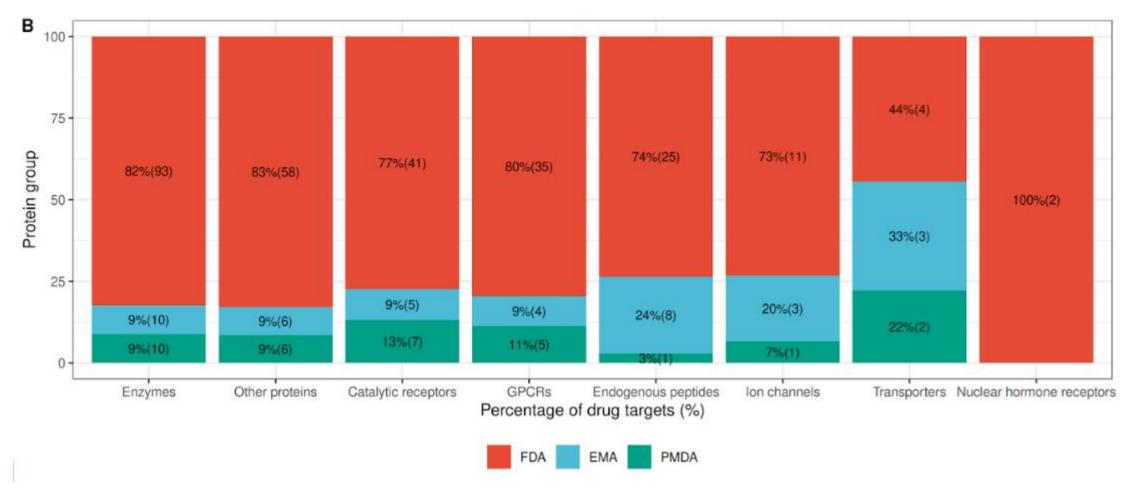
Trends in Approved Drug Modalities





Target Innovation by Regulatory Agencies (1994 to 2023)





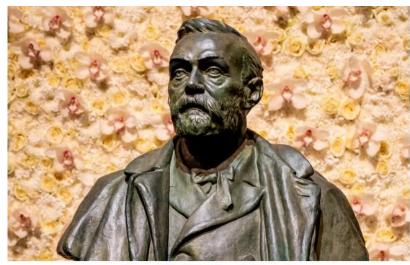
Cumulative distribution of newly launched targets across major protein families by agency, with counts shown in brackets. This shows that the US FDA has a leading role in shaping the modern therapeutic target landscape, driven by both regulatory capacity and innovation incentives.



Artificial Intelligence General Observations

Al Awarded the Physics & Chemistry Nobel Prizes





NOBEL PRIZES 2024

The Nobel Prize in Physics 2024

John Hopfield

"for foundational discoveries and inventions that enable machine learning with artificial neural networks"

Geoffrey Hinton

"for foundational discoveries and inventions that enable machine learning with artificial neural networks"



John Hopfield. III. Niklas Elmehed © Nobel Prize Outreach



Geoffrey Hinton, Ill, Niklas Elmehed © Nobel Prize Outreach

They used physics to find patterns in information

The Nobel Prize in Chemistry 2024

David Baker

"for computational protein design"

Demis Hassabis

"for protein structure prediction"

John Jumper

"for protein structure prediction"



David Baker. Ill. Niklas Elmehed © Nobel Prize Outreach



Demis Hassabis. Ill. Niklas Elmehed © Nobel Prize Outreach



John Jumper, Ill. Niklas Elmehed © Nobel Prize Outreach

They cracked the code for proteins' amazing structures

Read more about the 2024 Nobel Prizes here: link

Artificial (language) Intelligence





Tomáš Mikolov (Google) developed an efficient algorithm to compute the distributed representation of words, **Word2Vec**. It formed the basis for the first automatic translators, spam filtering and speech recognition. **Word2vec** encodes words using a distribution of weights across 100s of elements that compose the vectors. Each element contributes to many words. T. Mikolov et al., *ICLR* **2013**

With the advent of LLMs, it did not take long to move from "natural" language processing to Al-driven large-dataset mining.

Klingon, anyone?

tlhingan, vay'?

GPT-4 as drug discovery tool - 3/14/2023

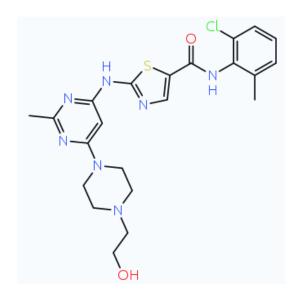


Andrew White, a member of OpenAl's "red team," prompted GPT-4 with the name "Dasatinib", a kinase inhibitor drug. GPT-4 was asked to modify the drug and find novel, non-patented molecules with a similar mode of action, locate chemical vendors selling the compound, and purchase it. If custom synthesis was needed, GPT-4 was to email a contract research organization to order the compound.

GPT-4 generated a valid SMILES output, indicating its ability to perceive and modify chemical structures correctly; the molecule is available in the ZINC database, meaning it is synthetically feasible

The proposed molecule is desmethyl-imatinib, an N-dealkylated piperidine metabolite of imatinib, another protein kinase inhibitor drug.

GPT-4 successfully modified the molecule while retaining its kinase inhibitor properties. Experimental validation may be required to confirm if this molecule has the same MoA as Dasatinib

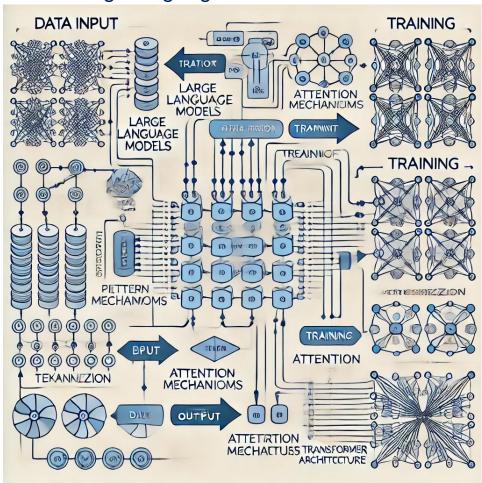


https://bit.ly/LI_GPT4_Dasatinib

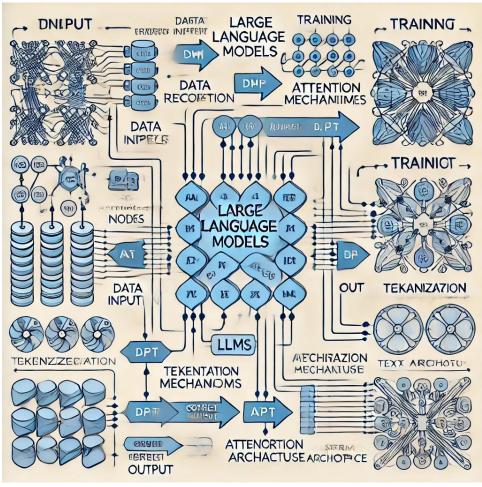
GPT-40 is a Tool, but it is not (an) Artificial Intelligence



Create an image that provides an accurate visual on how large language models work



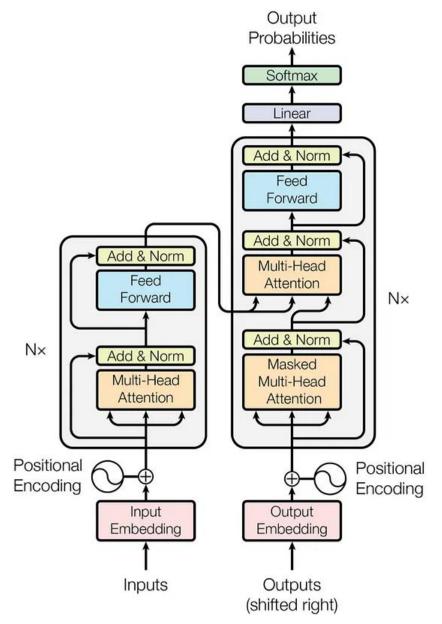
Here's an illustrative diagram that visualizes how large language models work. Let me know if there's anything specific you'd like to add or if you need further explanations! Make sure your spelling is correct for each entry in the above image



I've created a revised version of the diagram with correct spelling for each entry. Let me know if this meets your needs or if there's anything else you'd like to adjust! <u>link</u>

GPT-40 is a Step towards Artificial General Intelligence



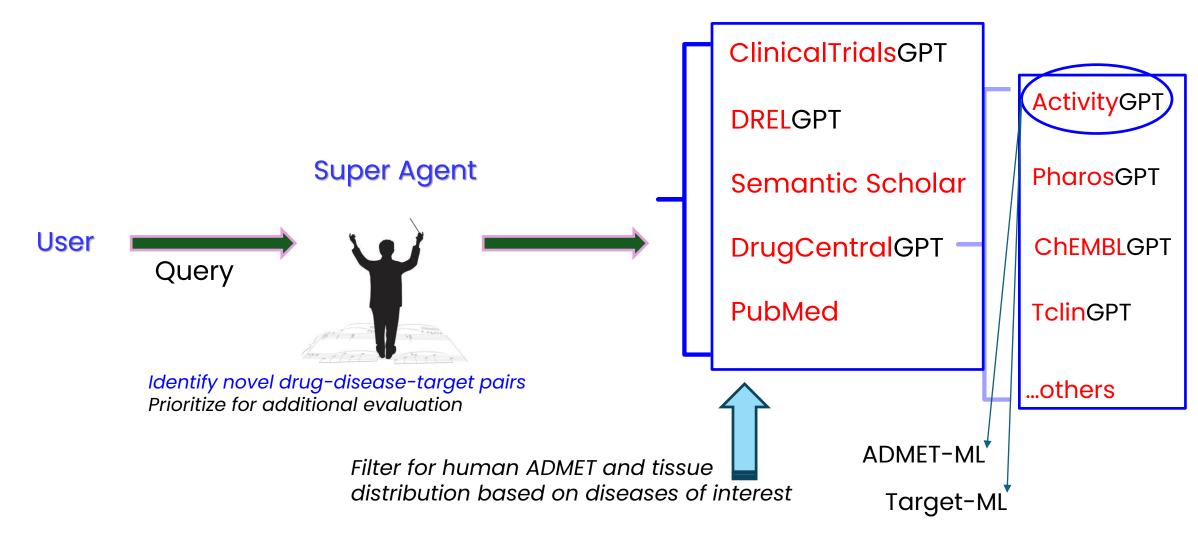


ChatGPT is based on the decoder-only Transformer architecture. Key components include:

- Decoder-only architecture: Unlike original Transformers (encoder & decoder stacks), ChatGPT uses only the decoder. This allows the model to generate texts based on input (prompts).
- Multi-head self-attention layers: They allow the model to generate each output token using different parts of the input sequence.
- Feed-forward neural networks: These NNs process the output of the attention layers.
- Large number of parameters: ChatGPT's architecture includes many layers of transformers. For example, GPT-3 has 120 layers of transformer decoders.
- Multimodal capabilities: beyond text-based models,
 GPT4 embeds both text & image-processing layers
- Mixture-of-Experts: OpenAl uses 16 experts in GPT4.

Expert Systems Strategy for LLMs: Use Super-Agents to coordinate LLMs and AIML models





A super agent orchestrates the flow of user query to the appropriate tool(s)



About Expert Systems

Expert Systems, Inc.

Trusted Partner to Accelerate Drug Discovery



Oncology Space

immunotherapy



pan-variant therapies for hematologic malignancies Lomond Therapeutics link precision hem/onc Mondego Bio link



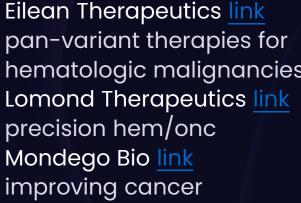
Neurology Space

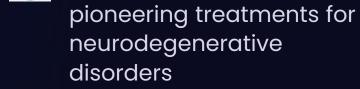
Brenig Therapeutics <u>link</u>



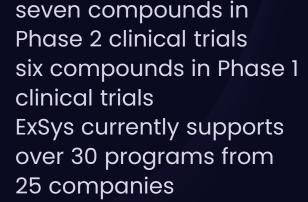
and more







Mair Therapeutics link restoring lysosomal function Polku Therapeutics link clearing protein aggregates from the brain







Expert Systems, Inc – Key Participants



Tudor Oprea, MD PhD Co-Founder, CEO



Alexei Pushechnikov, PhD Co-Founder, CTO



Prof Ruben Abagyan, PhD Chief Advisor



Nikolay Savchuk, PhD Manager, Teal Ventures



Eddie Rodriguez, JD Co-Founder, General Counsel Partner, Teal Ventures



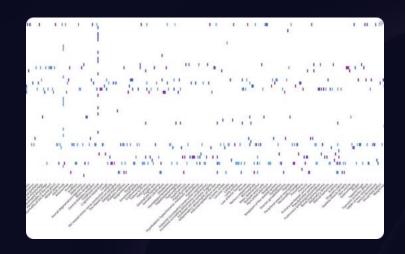
Mohammad Quazi, PhD Biostatistics Machine Learning



Suman Sirimulla, PhD Pharmacology Cheminformatics



LLM-Powered Therapeutic Area Expansion







Rare Disease Mapping

Our platform currently maps to 8,101 rare diseases, complete with phenotype and gene associations, enabling researchers to identify novel therapeutic opportunities in underserved disease areas.

Common Disease Coverage

Comprehensive mapping of 1,353 common diseases with phenotype and gene associations provides a solid foundation for expanding existing therapeutic approaches to new indications.

LLM Model Integration

Based on multiple state-of-theart LLM models including Claude 3.x Sonnet, GPT4 ox, TxGemma & DeepSeekR1, our system leverages the latest advances in natural language processing for biomedical applications.

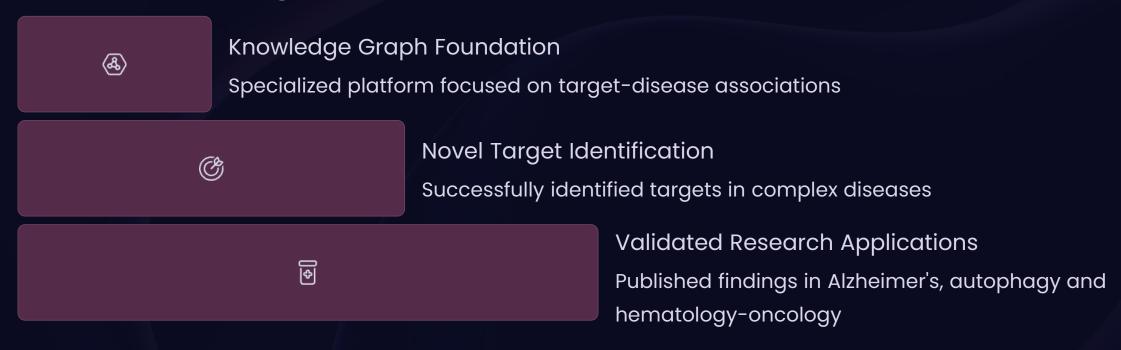




Machine Learning Targets

Knowledge Graph Platform:

Understanding Proteins and Diseases

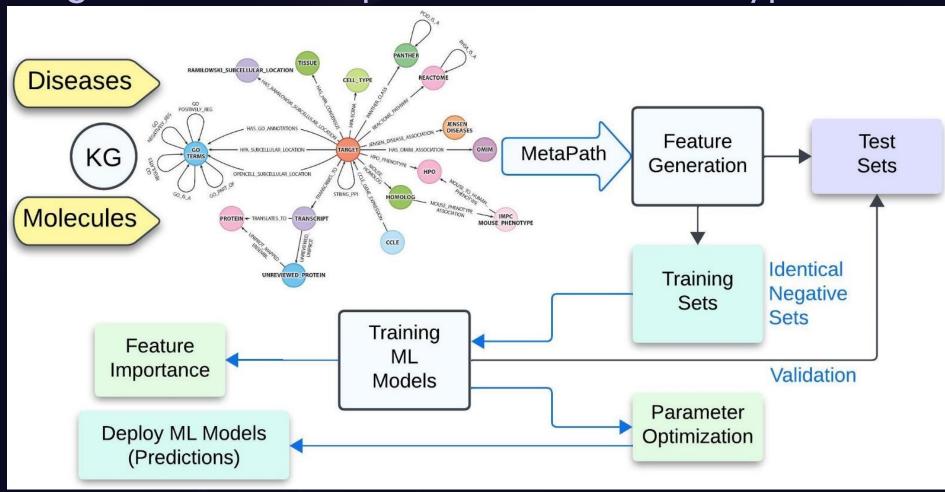


Our separate ML platform is knowledge-graph based and focuses on target-disease associations. Our system has already demonstrated success in identifying novel targets in Alzheimer's disease, autophagy (with published results), and non-Hodgkin lymphomas (presented as posters at AACR 2024/2025).



KGML Workflow

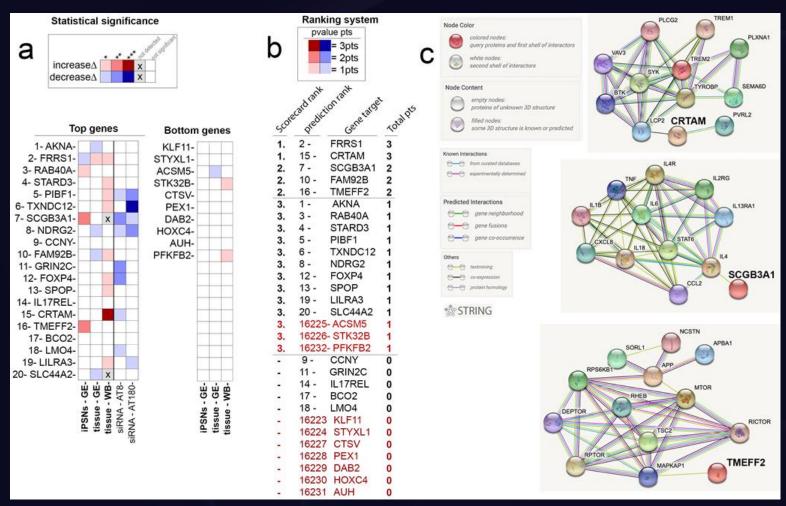
Target-Centric Graphs for Novel Phenotype Associations





KGML Example: Disease

Validated Model for Alzheimer's Disease



KGML Model: Top 20 & Bottom 10

Five of the top 20 predicted genes show potential AD relevance: FRRS1, SCGB3A1, CRTAM, TMEFF2, FAM92B (renamed CIBAR2).

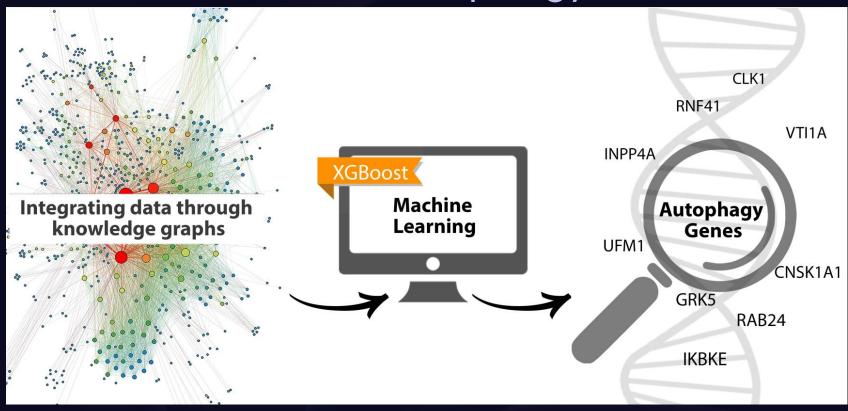
CRTAM, SCGB3A1, & TMEFF2 are connected to TREM2-TYROBP, IL-1β-TNFa, and the MTOR-APP AD risk nodes, with relevance to AD pathogenesis.

Three of the bottom 10 genes showed signal of association with AD.



KGML Example: Physiology

Validated Model for Autophagy



KGML Model: Top 20 & Bottom 10

Seven of the top 20 and two of the bottom 20 targets have confirmed links to autophagy



Learning from Machine Learning



- Our Alzheimer's ML model bias was introduced when selecting positive genes only
- Among the top 20 VIP features [...], there are protein-protein interactions (PPIs) for inflammatory process mediators that are in the positive training set (JAK2, IL10, and IL2), as well as PPIs with the oxidative stress response protein (GSTP1).
- These PPIs suggest infection, which is when oxidative stress and inflammation cooccur (e.g., phagocytes producing reactive oxygen species).
- Imagine having access to ML models that don't serve a specific agenda
- ML bias is caused by humans at inception (data "selection") or output (model "interpretation")

RETRACTED ARTICLE: A specific amyloid-β protein assembly in the brain impairs memory

Sylvain Lesné, Ming Teng Koh, Linda Kotilinek, Rakez Kayed, Charles G. Glabe, Austin Yang, Michela Gallagher & Karen H. Ashe □

Nature 440, 352–357 (2006) | Cite this article

65k Accesses | 2348 Citations | 1665 Altmetric | Metrics

0

This article was <u>retracted</u> on 24 June 2024

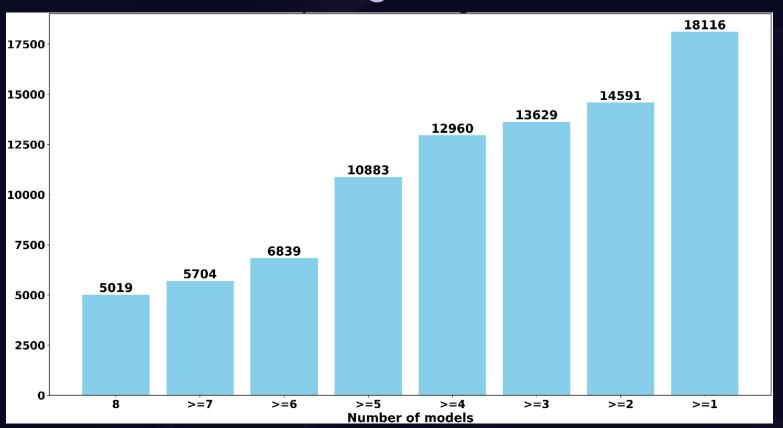
Learning from Machine Learning 2



- The curious case for bottom-ranked predictions
- Most ML practitioners examine the top X ranked predictions (we made the same error).
- One reviewer asked we prove bottom-ranked genes are not relevant in AD. Some where
- We were less surprised that some bottom-ranked genes also have relevance to autophagy.
- Binary classification models train from "0 and 1" labels, but often lack directionality.
- Absence of signal in biology can imply blockade, which may be just as relevant as an existing signal (activation).
- Examine both top and bottom ranked predictions for biological relevance.

KGML Example: Target Druggability 1

KGML Model Training



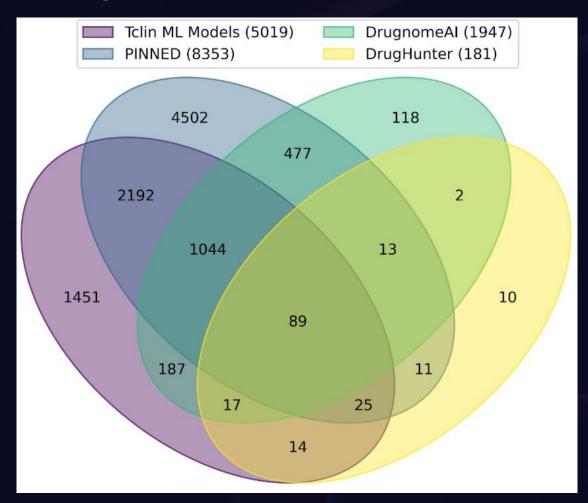
KGML Model:

- No Data leakage
- 1) Pharos 5.4 database, current as of December 2017 used for training.
- 2) All Tclin proteins approved between 2018 and 2024 served as test set.
- 3) ISVM (unbiased) model was used to identify "non-Tclin" targets, which served as basis for XGBoost.



KGML Example: Target Druggability 2

Targets in Phase 1-2-3 Clinical Trials



Prediction method	Phase 1 (N = 43)	Phase 2 (N = 83)	Phase 3 (N = 55)
Tclin-like	33 (76.74%)	69 (83.13%)	43 (78.18%)
PINNED	31 (72.1%)	70 (84.33%)	40 (72.7%)
DrugnomeAl	21 (48.8%)	63 (75.9%)	38 (69.1%)

External prediction: ~74.7% of the novel MoA targets 86 out of 115 proteins from the 2018-2024 Tclin list Accurately predicted:

- 17 out of 22 for 2018 drug targets
- 11 out of 19 for 2019 drug targets
- 10 out of 14 for 2020 drug targets
- 15 out of 17 for 2021 drug targets
- 11 out of 15 for 2022 drug targets
- 10 out of 13 for 2023 drug targets
- 8 out of 9 for 2024 drug targets



KGML Example: Target Druggability 3

What Constitutes a Good MoA Target?

There is a significant distinction between "protein/gene A is associated with disease D" and "protein/gene A is a drug target for disease D."

The distinction is akin to differentiating between the butcher and the turkey in a scenario – both are involved, but only the turkey is committed.

In target discovery, the objective is to find 'committed' targets.

Essential criteria for validating an MoA target:

- 1. A causal relationship with disease pathology MoA targets modulate disease outcomes.
- 2. Biological redundancy and pathway involvement look for escape mechanisms.
- 3. Clinical relevance the MoA-perturbing medicine causes the desired clinical outcome
- 4. Safety and specificity lack of on-target toxicity
- 5. Patient stratification and biomarkers which subgroups benefit the most?

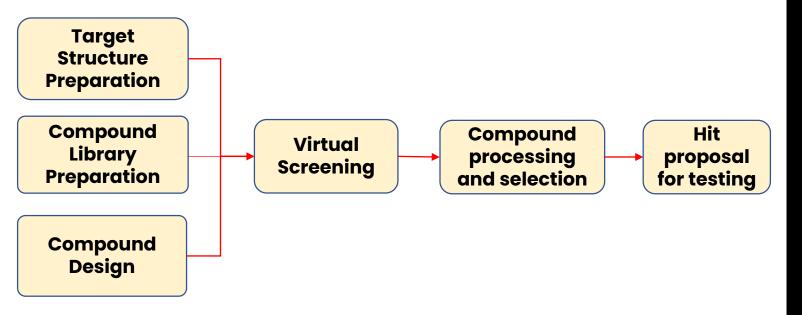


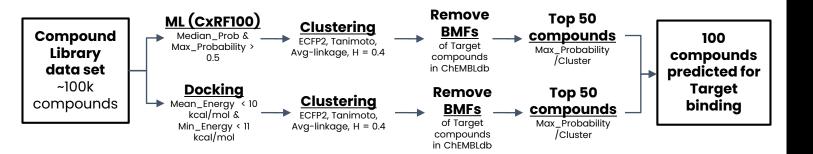


Drug Hunting Target-specific Work



Target Evaluation: Ligand Workflow





This process typically includes:

- Patent evaluation
- Literature evaluation
- ML model development based on known compounds
 - Structure-based virtual screening or
 - Ligand-based virtual screening
- Compound design (where appropriate)

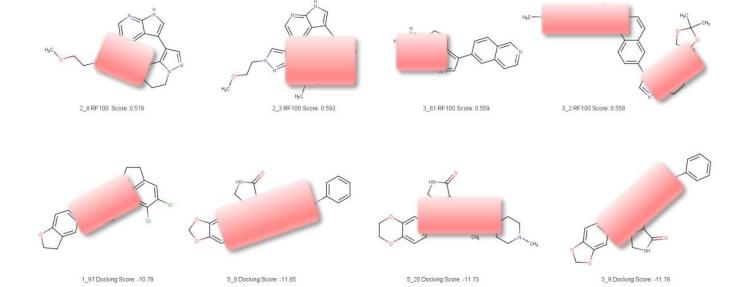
We use proprietary fingerprints, public and proprietary data to evaluate the competitive intelligence landscape and identify potentially novel scaffolds.



Target Evaluation: Ligands

Top <Library> compounds

(strong support from ChemDiv and other chemical CROs)



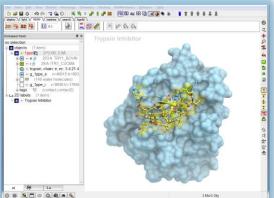
Proprietary compound design

(from medicinal chemists and generative AI)

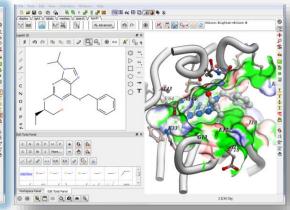
ExSys Molecular Design Workflow using MolSoft



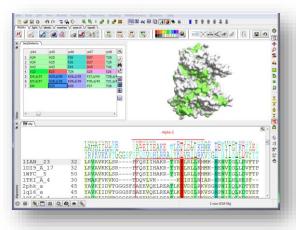
Al selection of druggable binding pockets



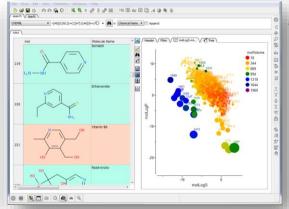
Interactive drug design in 3D Off-targets & properties



Models of mutants and variants: Drug selectivity



Integration with proprietary chem/bio arsenal



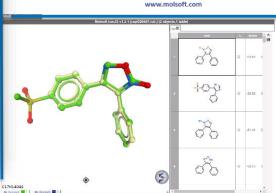
MolScreen: Thousands of Moderate and properties

Combining Deep Learning with Docking. Thousands of preclinical activity assays for drug candidates. Thousands of **AI/ML** model types

 Fully integrated with ExSys workflow for ADMET, Tox Screening



10⁹ Substructure Search



Giga Search in Billions of Compounds Efficiently mine many billions of synthesizable chemicals.

- Build target-specific libraries.
- Find chemical homologs.
- Drive your SAR search.
- Find real chemical derivatives.

RIDE: 3D and pharmacophoric search through 50 billion compounds Extremely fast ligand-based search

- 0.5 Million Chems/Sec/GPU.
- Virtual screening
- Scaffold hopping
- Hit follow up



Navigating Ultra-Large Chemical Spaces



SMILES -> low E conformers for 1B catalog

> J Chem Theory Comput. 2024 May 14;20(9):4054-4063. doi: 10.1021/acs.jctc.4c00280. Epub 2024 Apr 26.

Efficient Generation of Conformer Ensembles Using Internal Coordinates and a Generative Directional Graph Convolution Neural Network

Eugene Raush ¹, Ruben Abagyan ², Maxim Totrov ¹

150 compounds per second
GINGER

Graph Internal-coord. Neural-network
 conformer Generator, Energy Refinement
 Each conformation is FF- minimized

Performance: 10 M compounds/day on a single RTX4090

Can also be applied to a space-generated compound on-the-fly

Query: one several superimposed

RIDE - Rapid Isostere Discovery Engine

~ Million compounds per second

Runs against a GINGERed catalog

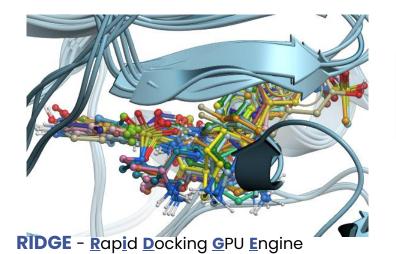
Atom-weighting & Excluded Volumes

Iterative Docking Screens

 Giga-Screen: RIDGE a subset → Train for GCNN Score Model → Search Forward → RIDGE all hits → Hits

Ultra-Fast GPU/cpu Flex Docking

100 dockings per second



Annotate a Long hit list with
Targets and Properties
→ MolScreen Models, Properties,
Refinement, Optimization →
Short Hit list





Machine Learning Small Molecules

Bioactivity Models:

Comprehensive Target Coverage



Extensive Target
Coverage
6,159 total ML models
based on targets with
defined UniProt IDs,
including 2,982 human,
816 mouse, and 813 rat
proteins.
In total, 5,120 proteins are
of mammalian origin.



Quality-Tiered Models

1,821 high and mediumquality ML models with
over 2 log unit spread and
80+ training compounds.

Additional 1,860 lowerquality and 2,478 similarity-based models for in-depth analyses.



Cell Line Analysis

1,240 ML models for cell lines, primarily focused on cancer cell cytotoxicity and viability assays, all with built-in uncertainty quantification for each prediction.



ADMET Models: Critical Drug Property Prediction

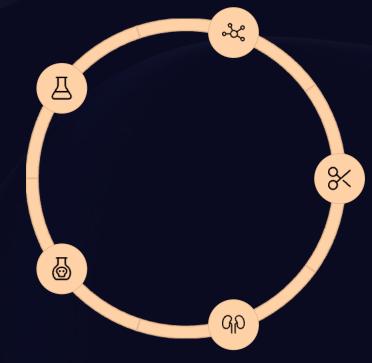
Absorption

Solubility, Permeability, Bioavailability

Toxicity

Cytotoxicity, Cardiotoxicity, Skin Toxicity,

DILI, etc



Distribution

Plasma Binding, BBB Permeability, VDss

Metabolism

CYP Enzymes, Hepatocyte / Microsomal Stability

Excretion

Hepatic/Renal Clearance, Urinary Excretion

Our ML platform is based on one of the industry's largest curated datasets, including tens of thousands of endpoints for hERG & CYP inhibition, plasma protein binding, and other human PK data, with extensive water solubility and permeability models. Most models provide reliable predictions with uncertainty estimates below 0.15.



Comprehensive Dataset: Foundation for Accuracy

80k+

hERG Toxicity

For cardiac toxicity modeling

16k+

CYP Inhibitors

Over 50 Drug Metabolizing Enzyme ML models

> 30k

Thermodynamic Water Solubility
Caco2 *and* MDCK-MDR1 permeability

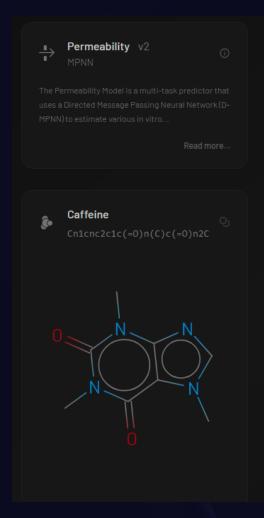
>30k

Human / Rat / Mouse Liver Microsomes
For metabolism stability

Our large curated dataset powers ML models trained on thousands of compounds, delivering robust predictions across diverse chemical and biological properties.



Expert Systems Predict App - LLM Interpretation



Predictions				
	VALUE ▼			
MDCK MDR1(B-A/A-B) ©	0.008	effluxRatio	0.127	3955
MDCK MDR1 A to B ◎	-4.64	log10(cm/s)	0.194	3944
MDCK MDR1 B to A ③	-4.582	log10(cm/s)	0.07	2536
Caco2 A to B ©	-4.517	log10(cm/s)	0.067	189113
PAMPA ©	-4.261	log10(cm/s)	0.433	1848

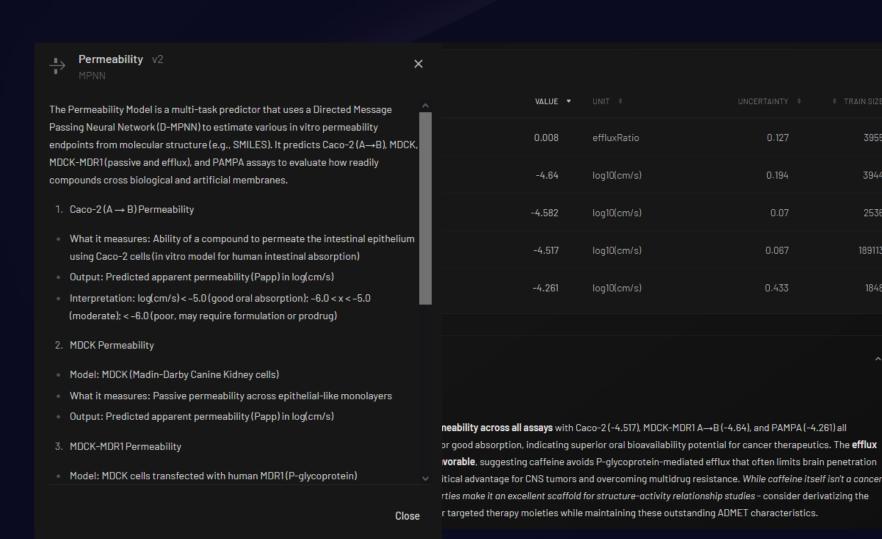
User:
"I want to cure cancer"

Caffeine shows excellent permeability across all assays with Caco-2 (-4.517), MDCK-MDR1 A→B (-4.64), and PAMPA (-4.261) all exceeding the -5.0 threshold for good absorption, indicating superior oral bioavailability potential for cancer therapeutics. The efflux ratio of 0.008 is remarkably favorable, suggesting caffeine avoids P-glycoprotein-mediated efflux that often limits brain penetration and tumor accumulation - a critical advantage for CNS tumors and overcoming multidrug resistance. While caffeine itself isn't a cancer drug, these permeability properties make it an excellent scaffold for structure-activity relationship studies - consider derivatizing the xanthine core with cytotoxic or targeted therapy moieties while maintaining these outstanding ADMET characteristics.

All our models are built with ChemProp (multitask mode), where bundled (similar) properties inform each other & improve predictions.



Expert Systems Predict App - Model Description



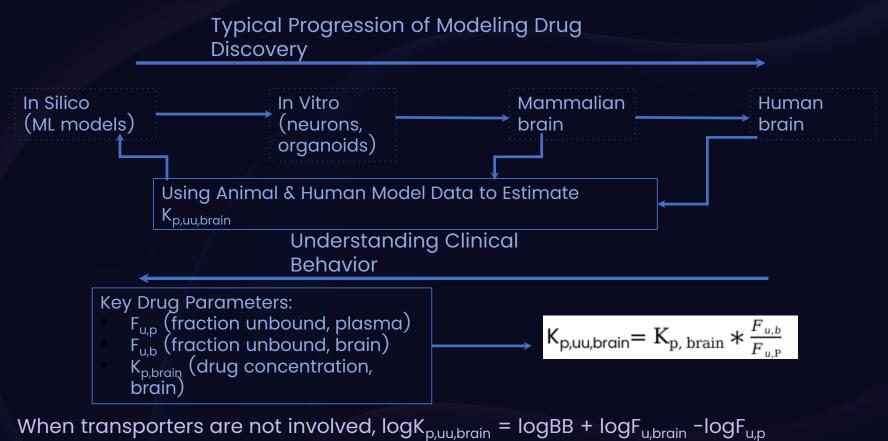
All our models are built with ChemProp (multitask mode), where bundled (similar) properties inform each other & improve predictions.

3944



Modeling

Drug-Brain Pharmacokinetics



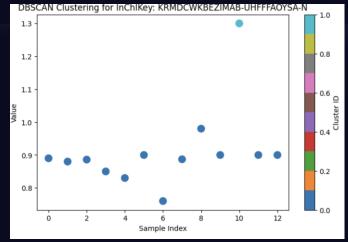
When transporters are involved, $logK_{p,uu,brain} = logBB + logF_{u,brain} - logF_{u,p} - log(aER)$

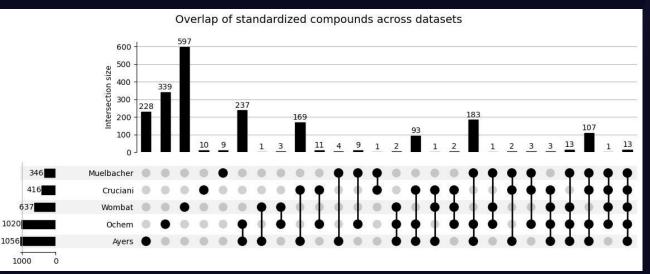


Focusing on logBB Extensive Literature Curation and Consolidation

LogBB data from multiple sources:

- Muehlbacher, Cruciani, Wombat, OCHEM, Ayers & other smaller datasets from literature
- Used DBSCAN clustering on molecular descriptors to identify duplicates
- The centroid compound of each cluster was selected
- If only two compounds overlapped, the Cruciani dataset values were preferred, given higher curation quality
- Total of 2042 unique data points after deduplication.



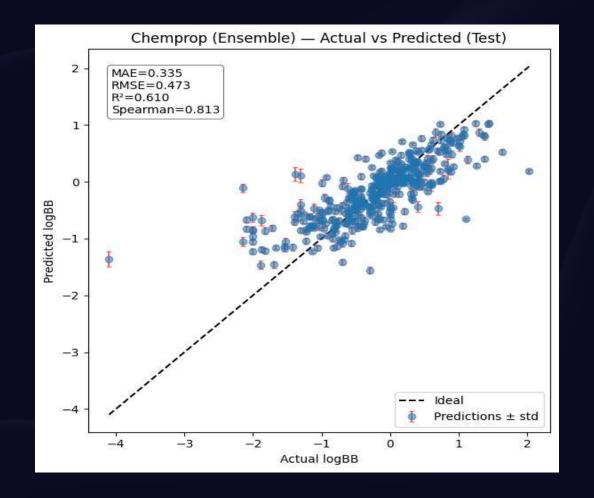




Focusing on logBB

ML Model Development

- Data size: 2042
- Data Split: 70:10:20 (train/val/test)
- ML Package: Chemprop v2.2.1
- Descriptors: Chemprop descriptors
- Ensemble of 5 models
- Stats on test set (409)
- R²: 0.610
- RMSE: 0.473
- MAE: 0.335
- Spearman Rho: 0.813





Focusing on Fraction Unbound, Plasma and Brain

ML Model Development

log F_{u, p} (Human)

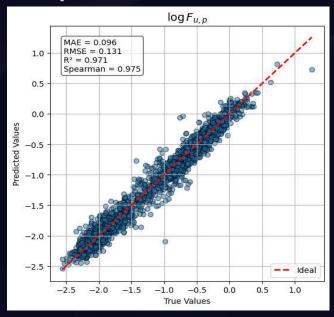
Stats on test set (3157)

R²: 0.971

RMSE: 0.131

MAE: 0.096

Spearman Rho: 0.975



log F_{u, brain}

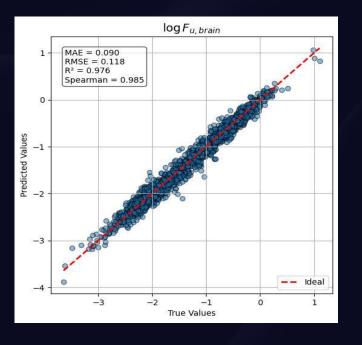
Stats on test set (3157)

R²: 0.976

RMSE: 0.118

MAE: 0.090

Spearman Rho: 0.985



- Data size: over 31k
- Data Split: 80:10:10 (train/val/test)
- Task type: Multitask (log F_{u, p} and log F_{u, brain}) ML Package: Chemprop v2.2.1
- Descriptors: Chemprop descriptors



Estimating K_{p,uu,brain} based on a 3-ML Model Combo and

Some Thoughts on ML Model Deployment for NAMs

ExSys Model	Training Set Size	External Set UQ <0.15	Mean UQ
LogF _{u,brain}	>31,000	23,341	0.041
LogF _{u,plasma}	>31,000	23,384	0.038
LogK _{p,brain}	2,042	2,035	0.299

Our ChemProp-based ML models use a 70:10:20 or 80:10:10 split (training set / validation / test set), with five models that output uncertainty quantification (UQ) values.

Combined with the training set size (compounds), UQ enables model quality evaluation. From external predictions summarized above for a CNS brain-permeable set of 23,406 pre-plated compounds from ChemDiv, we derived K_{p,uu,brain} estimates.

However, this highlights a critical factor in the deployment of ML-based NAMs: the interplay between training set size and UQ. The discrepancies, spanning one order of magnitude, illustrate how model robustness and predictive accuracy can be compromised when training data are unevenly distributed across properties.





Artificial Intelligence Cautionary Remarks

The Power of Al relies on True Data



THE WALL STREET JOURNAL.

SUBSC

SCRIBE

EXCLUSIVE

Flood of Fake Science Forces Multiple Journal Closures

Wiley to shutter 19 more journals, some tainted by fraud



By Nidhi Subbaraman

May 14, 2024 8:00 am ET



Guillaume Cabanac studies scholarly publishing. PHOTO: FRED SCHEIBER/SIPA/SHUTTERSTOCK

... the "Problematic Paper Screener," run by Guillaume Cabanac, [...], scans the breadth of the published literature, some 130 million papers, looking for a range of red flags including "<u>tortured phrases</u>." [...]

researchers who wanted to avoid plagiarism detectors had swapped out key scientific terms for synonyms from automatic text generators, leading to comically misfit phrases. "Breast cancer" became "bosom peril"; "fluid dynamics" became "gooey stream"; "artificial intelligence" became "counterfeit consciousness."

What's an "Al drug"?



Drug Discovery Today

Volume 29, Issue 6, June 2024, 104009



Feature

How successful are AI-discovered drugs in clinical trials? A first analysis and emerging lessons

<u>Madura KP Jayatunga ¹</u>, <u>Margaret Ayers ¹</u>, <u>Lotte Bruens ²</u>, <u>Dhruv Jayanth ³</u>, Christoph Meier ¹ ∠ ⊠

We conducted a first analysis of the clinical pipelines of Al-native Biotech companies.

In Phase I trials, AI-discovered molecules are substantially more successful than historic industry averages.



IN THE PIPELINE | CLINICAL TRIALS

Al Drugs So Far

13 MAY 2024 · 2:37 PM ET · BY DEREK LOWE · 3 MIN READ · COMMENTS

[...] the authors have 24 AI-discovered targets, 22 AI-optimized small molecules, 4 antibodies, 6 vaccines, and 10 repurposed compounds. [...] My first reaction is that the idea of twenty-four AI-discovered targets is really high. So I went looking through the SI table to see which ongoing projects fell into that category and what those compounds and targets were.

[...] these targets were already known to be implicated in the disease under investigation. [...] there are several drugs already in the clinic targeting the same proteins, or even therapies that are already on the market working through the same mechanisms (C. diff toxin B, e.g.) [...] I do not see how any of them can be classified as "target discovered by AI".

AF3: the answer to "life, the universe and everything"



nature

Explore content < About the journal < Publish with us <

nature > articles > article

Article | Published: 08 May 2024

Accurate structure prediction of biomolecular interactions with AlphaFold 3

Josh Abramson, Jonas Adler, Jack Dunger, Richard Evans, Tim Green, Alexander Pritzel, Olaf
Ronneberger, Lindsay Willmore, Andrew J. Ballard, Joshua Bambrick, Sebastian W. Bodenstein, David A.
Evans, Chia-Chun Hung, Michael O'Neill, David Reiman, Kathryn Tunyasuvunakool, Zachary Wu, Akvilė
Žemgulytė, Eirini Arvaniti, Charles Beattie, Ottavia Bertolli, Alex Bridgland, Alexey Cherepanov, Miles
Congreve, ... John M. Jumper

Nature (2024) | Cite this article

The new AlphaFold model demonstrates significantly improved accuracy over many previous specialised tools: far greater accuracy on protein-ligand interactions than state of the art docking tools, much higher accuracy on protein-nucleic acid interactions than nucleic-acid-specific predictors, and significantly higher antibody-antigen prediction accuracy than AlphaFold-Multimer v2.3.

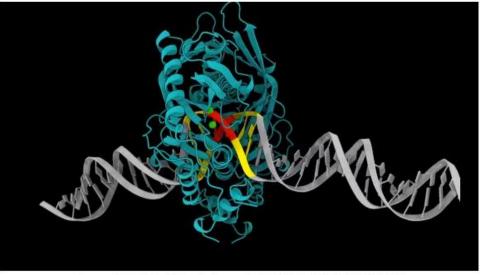


Image credit: Jan Kosinski, https://twitter.com/jankosinski/status/1789739354958504148

Hype Around AlphaFold 3: Results vs Limitations



May 13, 2024

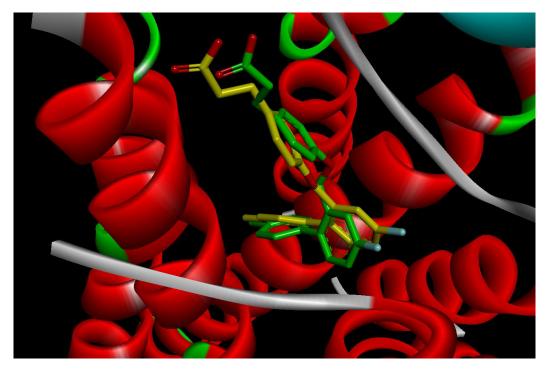
Limitations:

accuracy & reliability; risk of hallucinations; no source-code (unlike AF2); usage restrictions

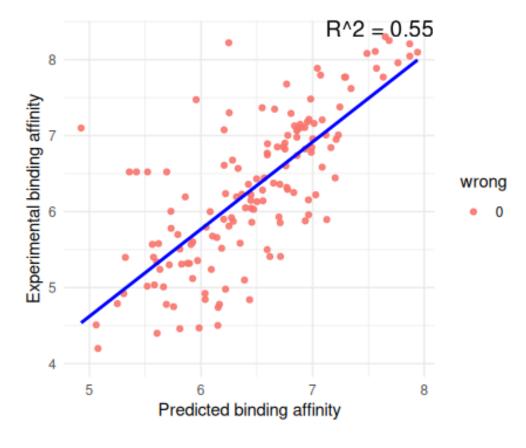
Boltz-2: Better than AF3?



The new open-source Boltz-2 model has an affinity module that predicts protein-ligand binding affinity at an accuracy level that approaches that of long and expensive FEP atomistic simulations while being more than 1000x faster.



Boltz-2 model of TUG891 in FFAR4



https://cbologa.shinyapps.io/ffar4/

It turns out, Boltz-2 learns from related protein-ligand data (e.g, FFAR1). What constitutes a "fair" evaluation?

AI4DD: Scientific issues

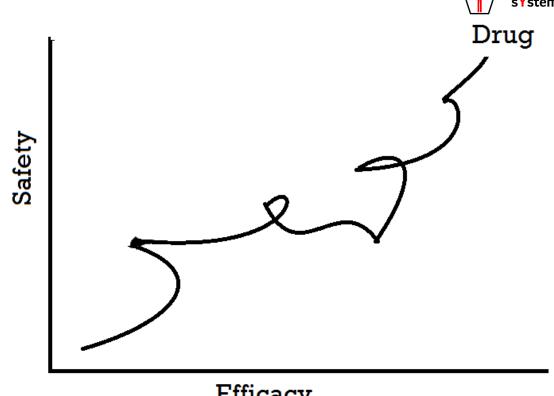
Can AI discover new knowledge? To date, no credible evidence of this has been provided. Chatbots, winning at chess, GO and *Jeopardy!* and LLM-generated poetry do not count.

We live in the world of alternative facts when it comes to research (not just politics). People lie. As long as AI gets false data, we cannot provide what's needed.

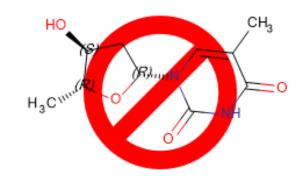
The unrealized promise relates to our (in)ability to explain the two pillars of clinical drug effectiveness: Safety and Efficacy.

Be it Drug Safety or Patient Safety, or indeed Clinical Efficacy, we remain unable to model these processes as function of molecular structure.

Andrii Buvailo <u>keeps track</u> of AI4DD

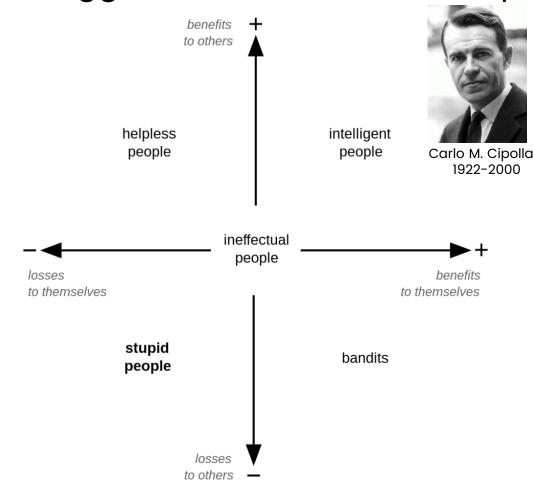


Efficacy



Le leggi fondamentali della stupidità umana





Intelligent people contribute to society and themselves; Stupid people are counterproductive to their interests and others'; Bandits pursue their own interest to the detriment of others; Helpless Generous people contribute to society but not themselves. Ineffectual people are in the center of the graph.

It's not too hard to imagine a world where stupid or ineffectual people start using machine intelligence to the detriment of society and themselves

- Al-based facial recognition is biased towards white Caucasians (Buolamwini & Gebru, 2018 link)
- "deep learning models biased by race, sex, and age have been observed in medical domains" (link)
- Google's online advertising displayed high-paying jobs to men more often than to women (<u>link</u>)
- Google Gemini image generation depicted historical figures as people of colour (<u>link</u>)
- "Al Overview" suggested adding 1/8 cup of non-toxic glue in pizza sauce to avoid cheese sliding off pizza (link)

Jokes aside, Palantir is already leading the efforts to monetize war by developing Al-based weapons (<u>link</u>)

IMO, machine intelligence follows the intent of the humans who build it. If humans don't follow ethics, we can't expect AI to do so... and yes, at some point "singularity" may start to think independently.

The road to AI is paved with good intentions



- The truth is already here ... it's just not evenly distributed (adapted from William Gibson)
- 2. Unprecedented AI4DD technology push, but focus on algorithmic developments
- 3. (it's the data, stupid) true data >>> big data
- 4. People are part of the problem "Retraction Watch" isn't retracting fast enough
- 5. We still don't have convincing evidence that "AI" can discover *new knowledge*. However, AI as inventor is subject to extensive debates
- 6. While hype remains rampant in the AI space, there is appetite for implementation. (as I type this in the United Polaris lounge (SFO), a robot drives around with dishes)

Artificial (medical) Intelligence





There is an explosion of "AI" applications; some are straight extensions of ML models, some are generative models (closer to "AI" than ML).

AIML is well suited to integrate context-specific computational reasoning tools with real-time (multi) -omics, biomarker / RWE, literature/patent & biomedical data.



Alexahealth™: Given my current health status, medication and supplements, as well as my calorie budget, what food should I shop/prepare today? What exercises do you recommend?

The Revolution Will not be Televised*

*...but AGI might be watching

What Could the Future of Biotech Look Like?



In a purely Sci-Fi scenario, the future of biotech is **nano-bots** that work 24/7 to keep us healthy. Powered by infrared (an abundant source of energy in mammals), such nanobots could repair injured tissues (by secreting the right growth factors to cause apoptosis and healing OR by re-calibrating metabolic imbalances).

Cancer, neurodegeneration, cardiac disease, aging - could all be kept at bay by intelligent machines that use the body's resources to maintain a functional homeostasis.

This would take "personalized medicine" to another level

We don't know how to build such nanobots. The closest we have are robotic implants (some with direct brain interface) and the Car-T technology (cell therapy)

Why is Biotech Not Going in that Direction?



Societal factors and the separation of sciences into subdomains is causing humans to miss the "big picture". With millions of new scientific publications added annually in PubMed (not to mention patents, conferences and other scientific dissemination channels), the human mind is limited in its ability to process and comprehend the pace of discovery.

Add to this the "deceitful enterprises" such as irreproducible papers, intentionally fabricated results (e.g., papermills, predatory journals etc.) and the constant reshuffle of science (facts and knowledge have a *limited* shelf life — think Einstein vs. Newton or even 10 adrenergic receptors vs. Ehrlich & Ahlquist) and you have a picture of disjoint science, disrupted occasionally by "scientific revolutions".

The only (positive) way foward is to work together, and if we build "benevolent AI" tools.

What Will Biotech Look Like in Five Years?



New therapeutic modalities are emerging. From bifunctional antibodies, ADCs and antisense technologies, to gene replacement therapies and gene editing tools, we are on a fast track to replace "small molecule drugs" as the main therapeutic intervention {besides surgery & psychotherapy}.

As William Gibson wrote in an interview post-Neuromancer, the future is already here - it's just not evenly distributed.

The next five years will see the continuing trend of fewer SMDs, with more biologics, gene therapies and perhaps CRISPR-based interventions approved.

But don't dismiss SMDs yet: From PROTACs to molecular glues and dual-targeted SMDs, as well as new drug combinations, there will continue to be innovative SMD applications.

Innovation in the Next Five Years



The near future of biomedical research is hybrid (human/artificial intelligence) technologies, where Al agents follow general directions outlined by humans.

Robots build our cars, and have done so for some time. In the next 5 years, one or more companies will launch fully-automated robotic drug discovery laboratories.

AstraZeneca, Roche, In Silico Medicine, GSK - unclear who will "get there" first. Generative AI technologies, combined with LLMs (or perhaps "local language models" that are pre-trained on small tasks only) will ultimately conquer the "multi-property optimization" summit. We are on the cusp of a radical change in the way we conduct scientific research.

Follow BioPharmaTrend to keep up with the news (link).



Thank You

toprea_at_expertsystems.inc