

A large, colorful molecular structure graphic on the left side of the page. It consists of numerous spheres in various colors (blue, green, red, yellow, orange, pink, white) connected by thin white lines, representing atoms and bonds. The structure is set against a light blue background with a faint silhouette of a human head in profile, suggesting the connection between molecular science and human health.

FROM  
MOLECULE TO  
PATIENT

ASCPT 2019  
ANNUAL MEETING



# An Integrative Deep Learning Approach for De Novo Drug Discovery

Joel Dudley PhD

Icahn School of Medicine at Mount Sinai

# Conflicts of Interest Statement

## **Scientific founder or co-founder**

- Onegevity Health
- OOVA
- Ontomics
- NuMedii

## **Scientific advisory board member**

- Ayasdi
- LAM Therapeutics
- Solve Bio
- Hoy Health
- Biotia

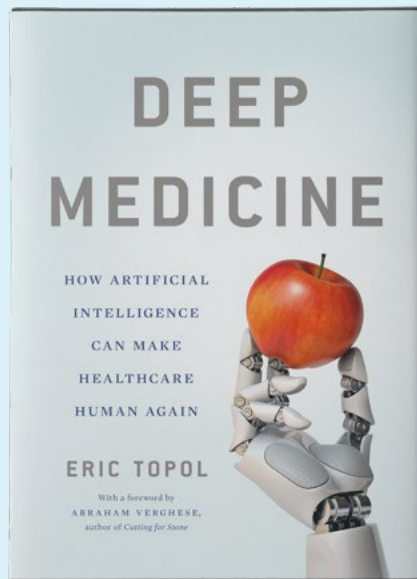
## **Past or present consultancy**

- Janssen Pharmaceuticals
- Thorne Research
- Allergan
- AstraZeneca
- LEO Pharma
- Speaking honoraria
- Celgene
- Illumina
- Roche
- Takeda
- Lundbeck



# Deep Learning hits the mainstream in biomedicine

FROM  
MOLECULE TO  
PATIENT



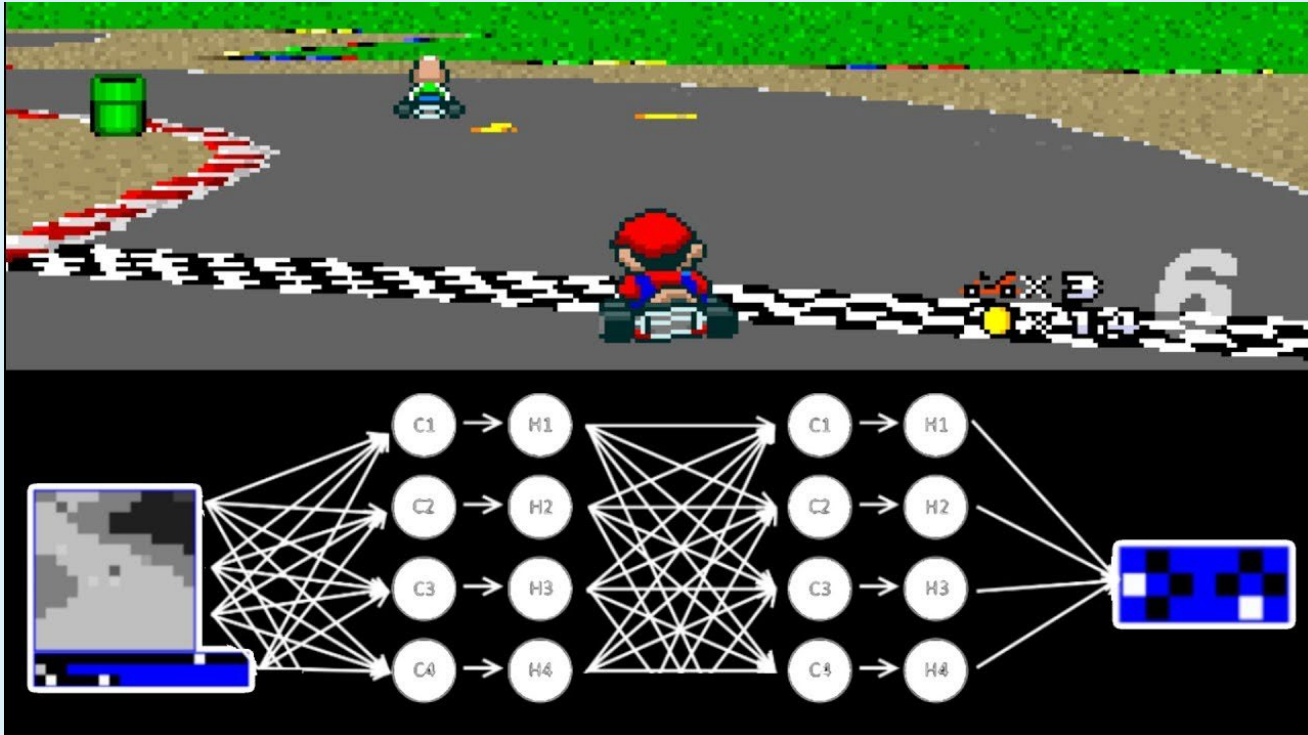
**Dr. Eric Topol**  
**Scripps Translational Research Institute**



# Deep Learning – what's the big deal?



# Deep Learning – what's the big deal?



# Deep Learning – what's the big deal?

## Google computer works out how to spot cats

**A Google research team has trained a network of 1,000 computers wired up like a brain to recognise cats.**

The team built a neural network, which mimics the working of a biological brain, that worked out how to spot pictures of cats in just three days.

The cat-spotting computer was created as part of a larger project to investigate machine learning.

Google is planning to use the learning system to help with its indexing systems and with language translation.



Millions of images were used to train the neural network

**Related Stories**

<http://www.bbc.com/news/technology-18595351>



# Deep Learning w/ EHR data



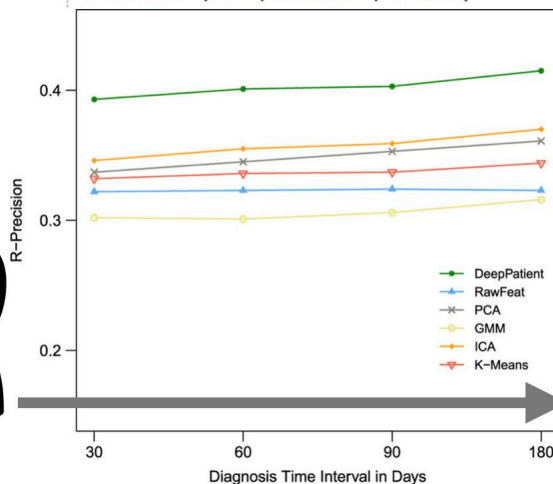
MIT  
Technology  
Review

SCIENTIFIC  
AMERICAN

## OPEN Deep Patient: An Unsupervised Representation to Predict the Future of Patients from the Electronic Health Records

Received: 28 January 2016  
Accepted: 27 April 2016  
Published: 17 May 2016

Riccardo Miotto<sup>1,2,3</sup>, Li Li<sup>1,2,3</sup>, Brian A. Kidd<sup>1,2,3</sup>, Joel T. Dudley<sup>1,2,3</sup>

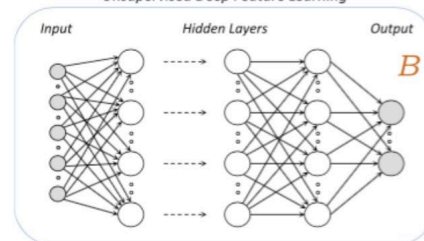


**Electronic Health Records**  
Clinical Notes  
Diagnoses  
Medications  
Laboratory Tests  
Demography  
Etc.

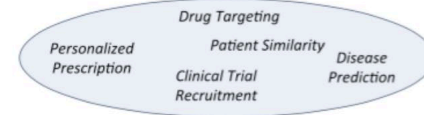
Raw Patient Dataset



Unsupervised Deep Feature Learning

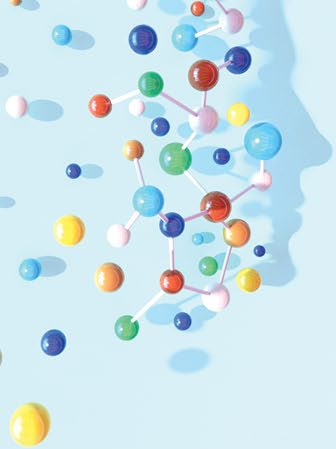


Deep Patient Dataset



Miotto R, Li L, Kidd BA, Dudley JT. Scientific Reports 6 (2016)





# Deep Learning w/ EHR data

## *Deep Logistic Regression Network*

Disease	AUC-ROC
Cancer of Liver	0.93
Regional Enteritis and Ulcerative Colitis	0.91
Type 2 Diabetes Mellitus	0.91
Congestive Heart Failure	0.90
Chronic Kidney Disease	0.89
Personality Disorders	0.89
Schizophrenia	0.88
Multiple Myeloma	0.87
Delirium and Dementia	0.85
Coronary Atherosclerosis	0.84

# Deep Learning and Donald Rumsfeld. Bet you didn't expect that.

*There are known knowns; there are things we know that we know.*

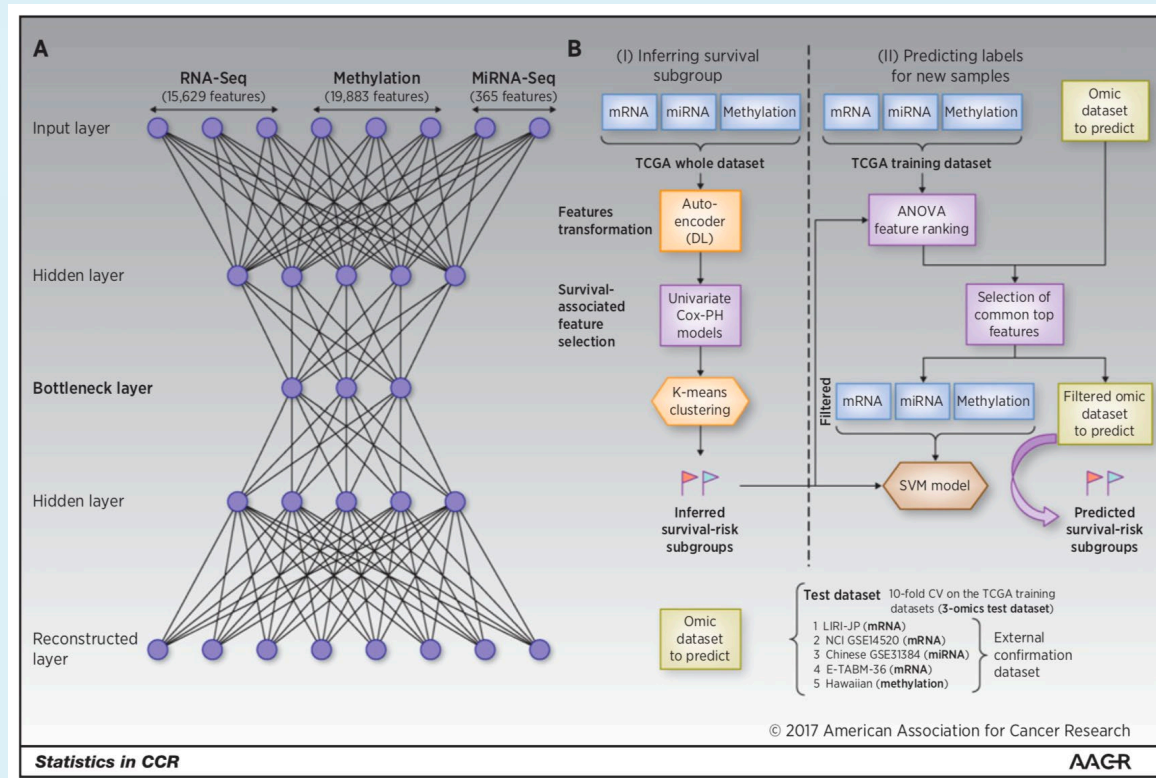
*There are known unknowns; that is to say, there are things that we now know we don't know.*

*But there are also unknown unknowns – there are things we do not know we don't know.*

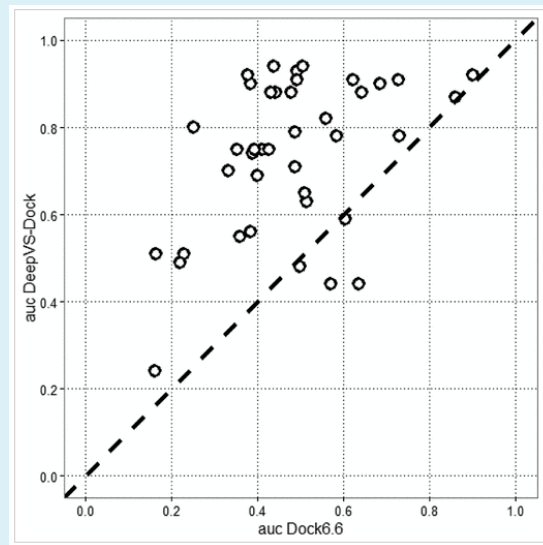
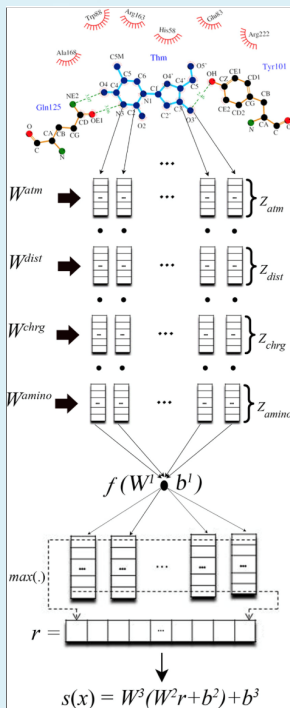
-Donald Rumsfeld



# Learning latent factors and compressed representations




# Boosting Docking-Based Virtual Screening with Deep Learning



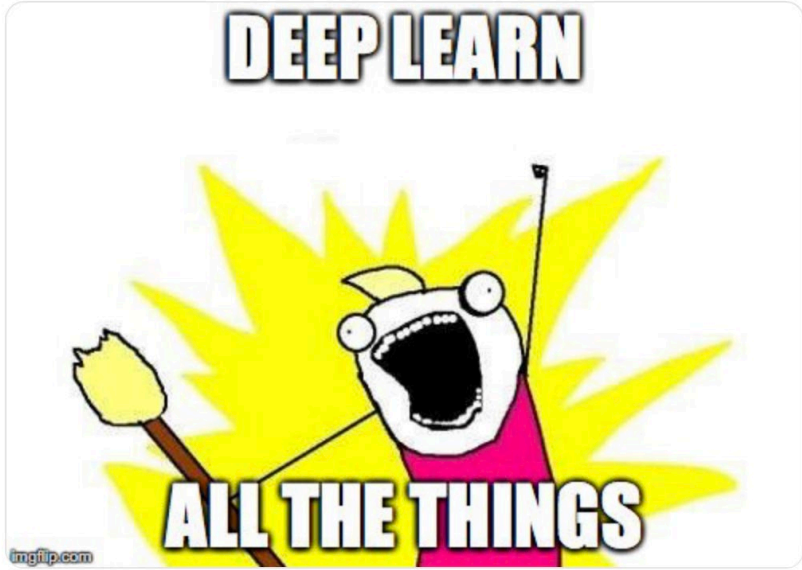
- DeepVS uses convolutional neural networks to learn abstract features (i.e. compound atom type, atomic partial charged and distance between atoms) that can discriminate active ligands
- The deep learning method outperforms two well-established virtual screening methods on a benchmarking dataset



## Prophecy and Deep Learning

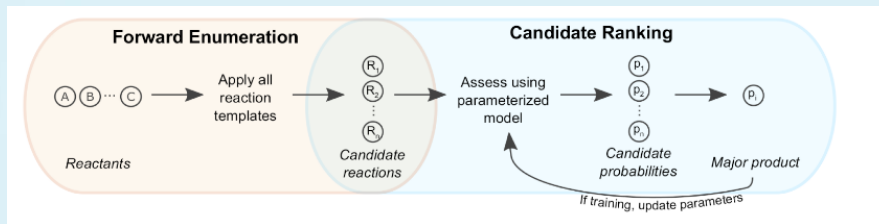
 **jdudley**  
@jdudley

This image summarizes all the bioinformatics papers that will come out in the next 6 months (I am among the guilty)

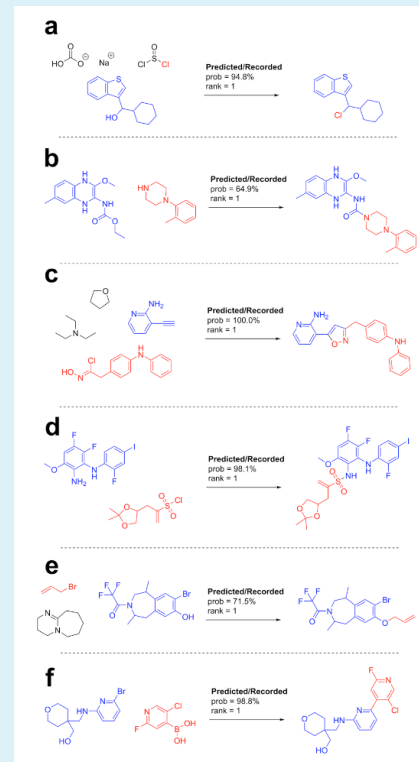


7:06 PM - 30 Mar 2016

# Prediction of organic reaction outcomes using machine learning

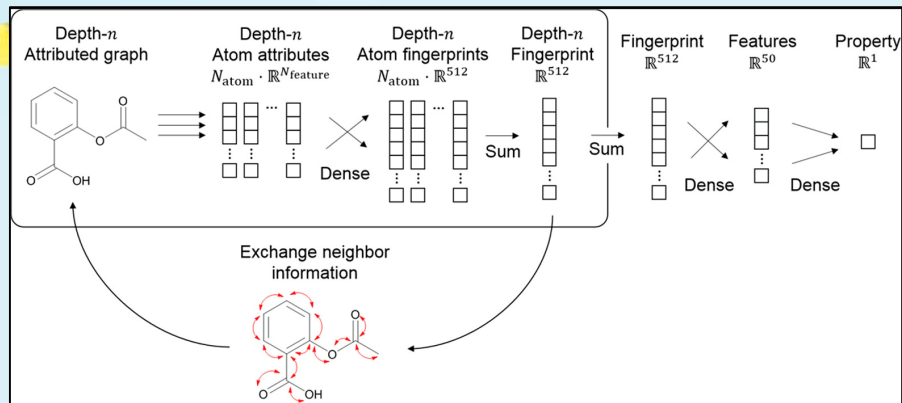


- Organic compound synthesis is difficult and requires extensive human planning.
- Existing computer-aided techniques essentially just iteratively apply known reaction templates
- Proposed reaction steps which work on paper often fail in the laboratory
- This manuscript uses published reactions from 15,000 US patents to generate lists of candidate products from reactants





# Convolutional Embedding of Attributed Molecular Graphs for Physical Property Prediction



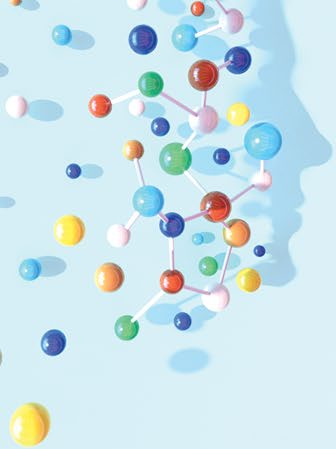
Structure	Z	# neighbors	# hydrogens	Formal charge	
	Node 1	6	1	3	0
	Node 2	6	2	2	0
	Node 3	8	1	1	0

Attributed Graph	Order	Aromatic	Conjugated	In ring	Connects
	Edge 1	Single	No	No	(1, 2)
	Edge 2	Single	No	No	No

$$M_{\text{ethanol}} = \begin{bmatrix} 6, 1, 3, 0, 0, 0, 0, 0, 0 & 6, 2, 2, 0, 1, 0, 0, 0, 1 & 0, 0, 0, 0, 0, 0, 0, 0, 0 \\ 6, 1, 3, 0, 1, 0, 0, 0, 1 & 6, 2, 2, 0, 0, 0, 0, 0, 0 & 8, 1, 1, 0, 1, 0, 0, 0, 1 \\ 0, 0, 0, 0, 0, 0, 0, 0, 0 & 6, 2, 2, 0, 1, 0, 0, 0, 1 & 8, 1, 1, 0, 0, 0, 0, 0, 0 \end{bmatrix}$$



# Hello, I don't exist!

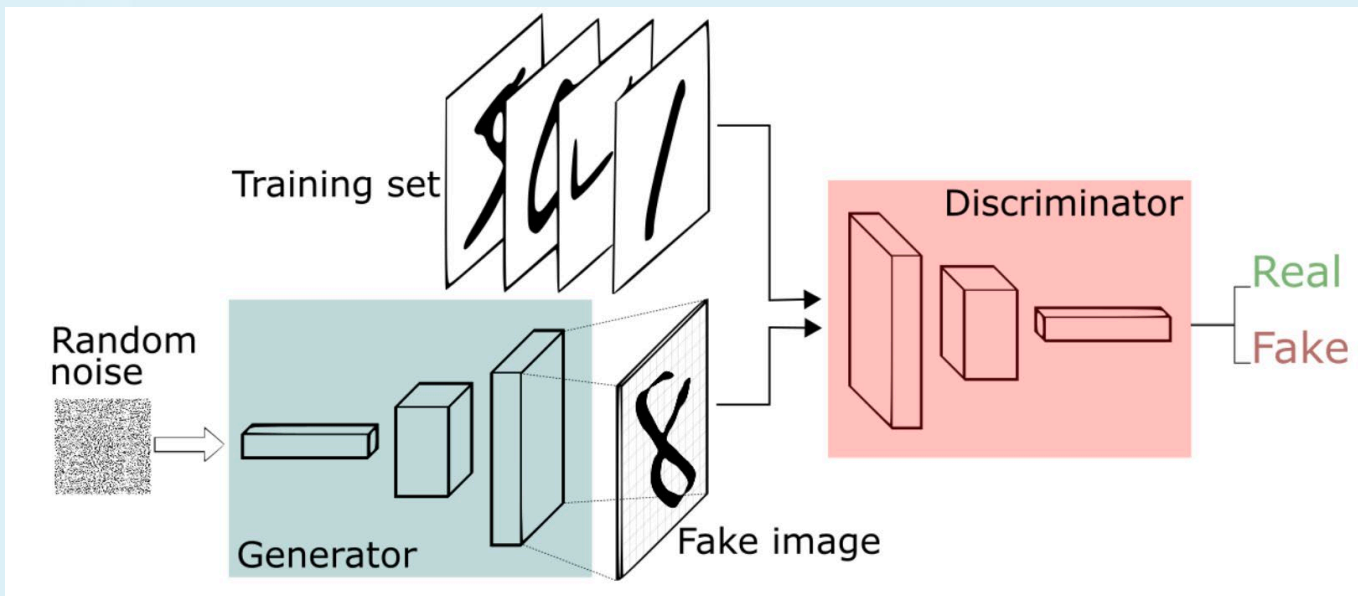
FR OM  
MOLECULE TO  
PATIENT



<https://thispersondoesnotexist.com/>

StyleGAN – Keras et al. 2018

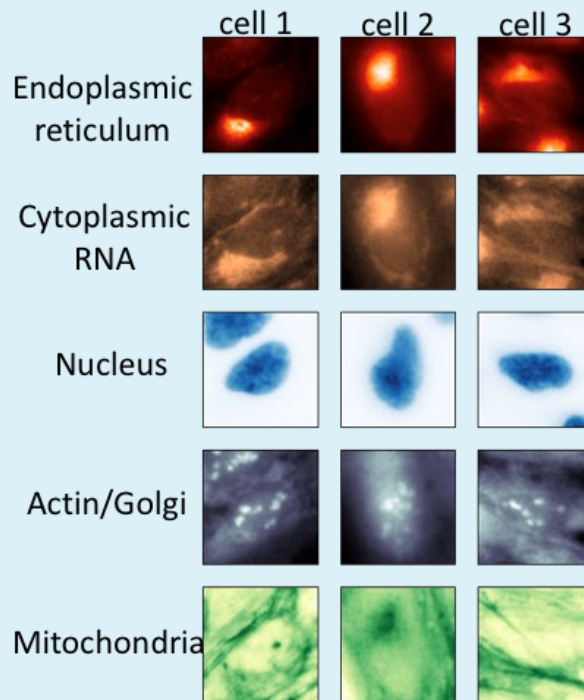
# Generative Adversarial Networks (GAN)



# GAN applied to phenotypic screening scenario

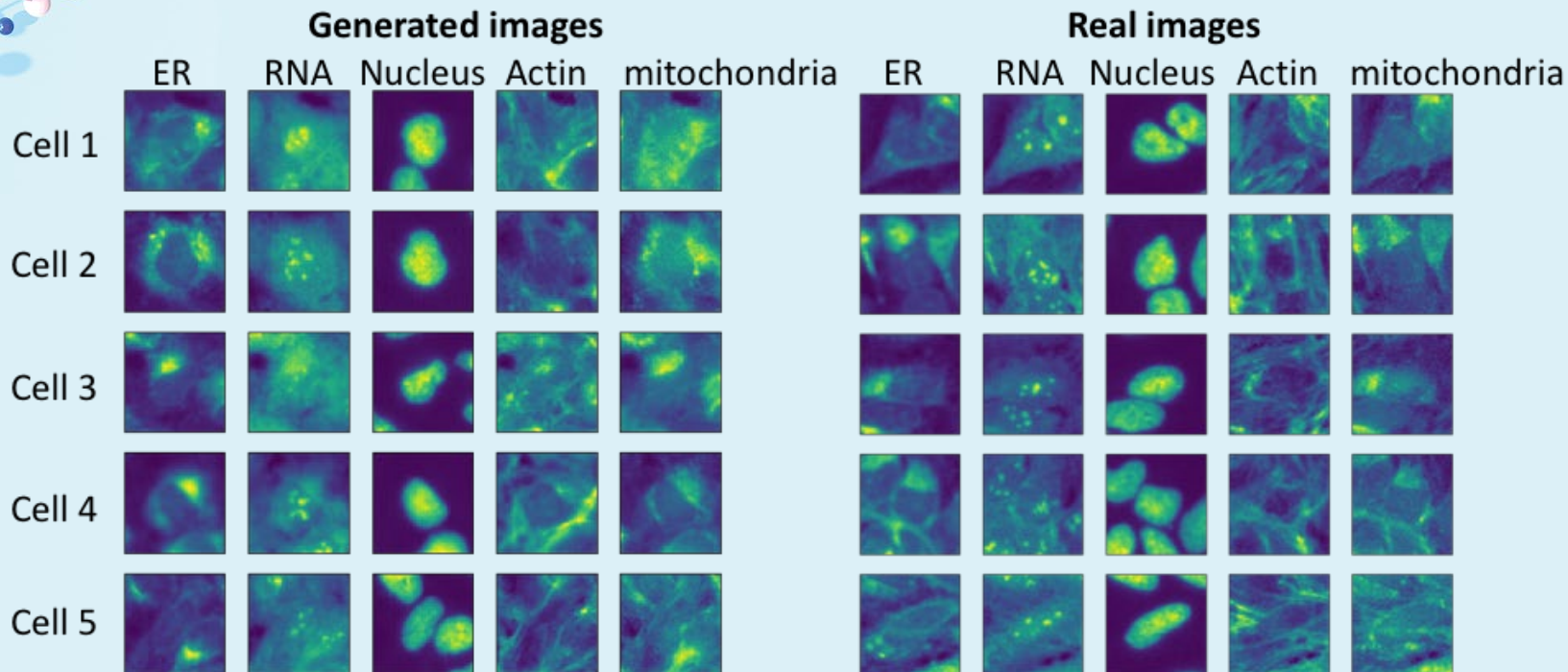
FROM  
MOLECULE TO  
PATIENT

- On average 4 repeats per treatment
- 6 images form each repeat
- 10~50 cells from 1 images
- 2 TB of image patches (64x64)

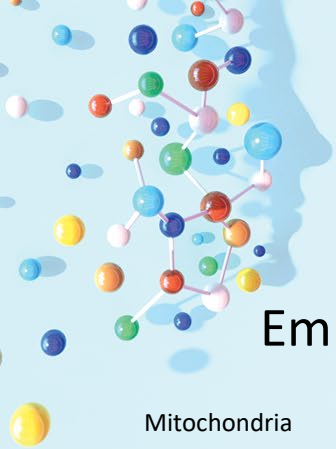




# GAN applied to phenotypic screening scenario



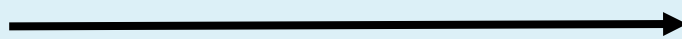




# Moving from one image to another in the latent space

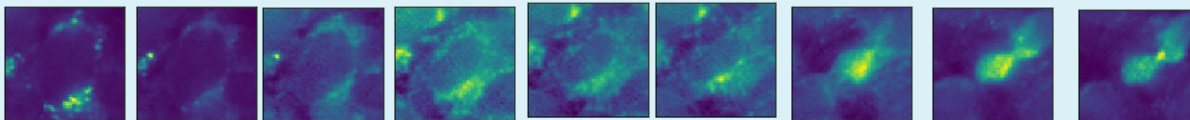
FROM MOLECULE TO PATIENT

Embedding 1

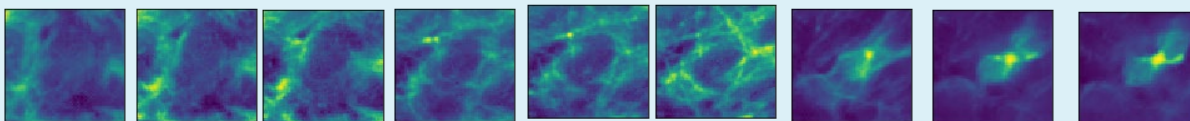


Embedding 2

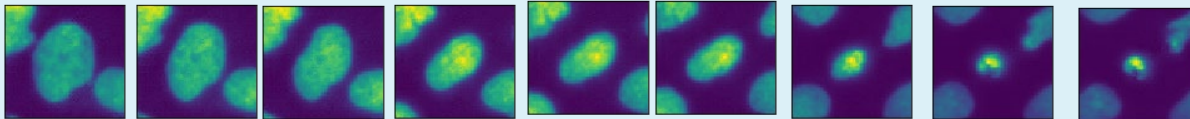
Mitochondria



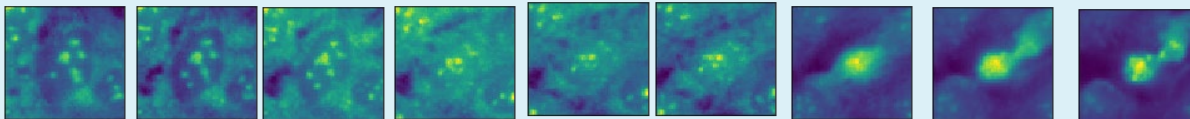
Actin/Golgi



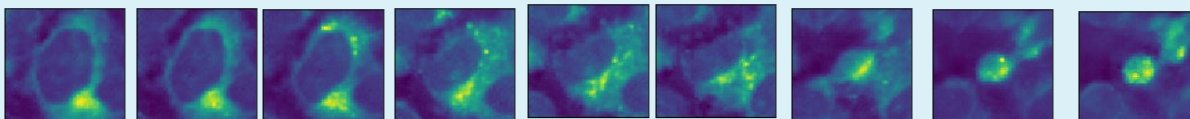
Nucleus



Cytoplasmic RNA

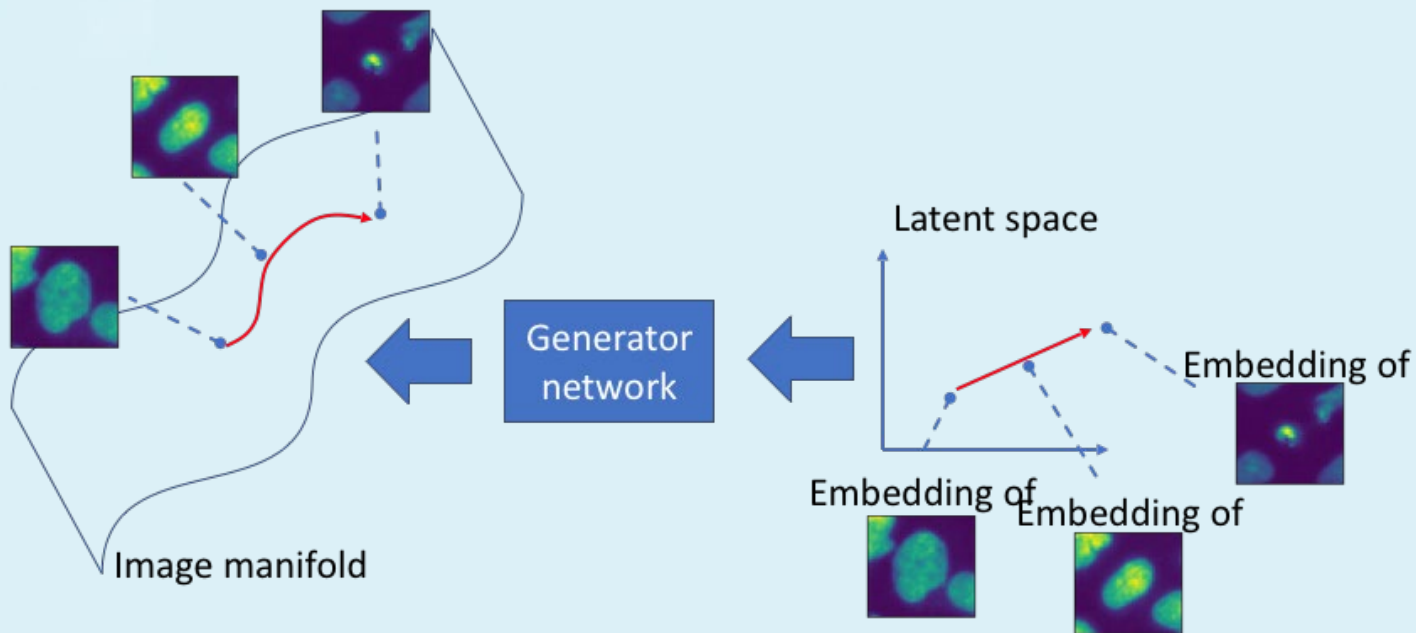


Endoplasmic reticulum

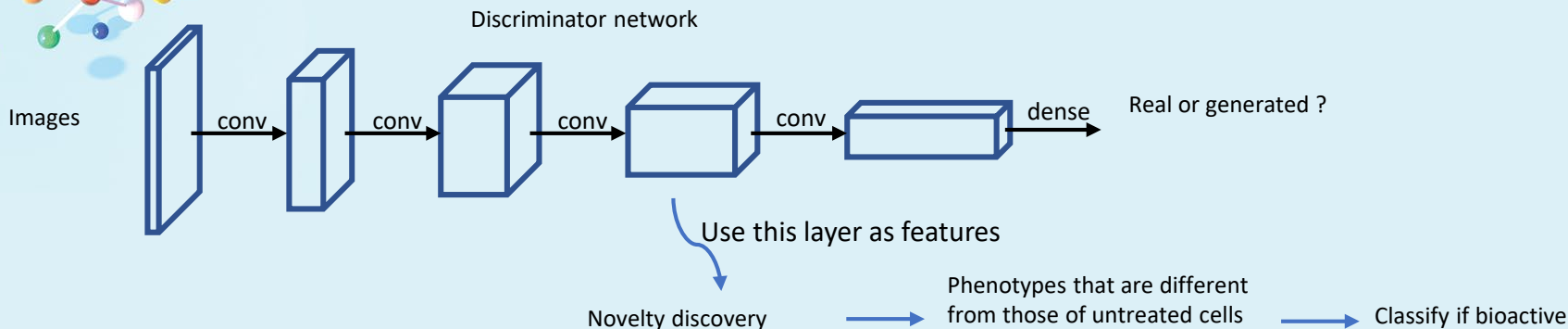




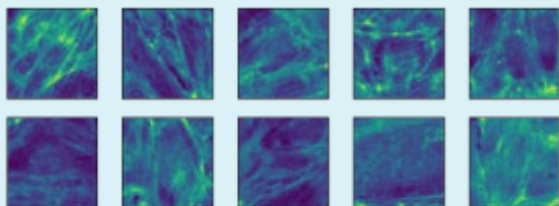
# GANs encode “action” on the images



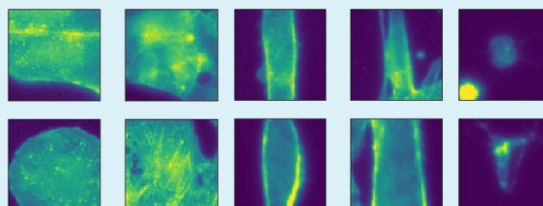
# Using the discriminator to discover bioactive chemicals



Untreated cells (actin filament)



Novel phenotypes (actin filament)

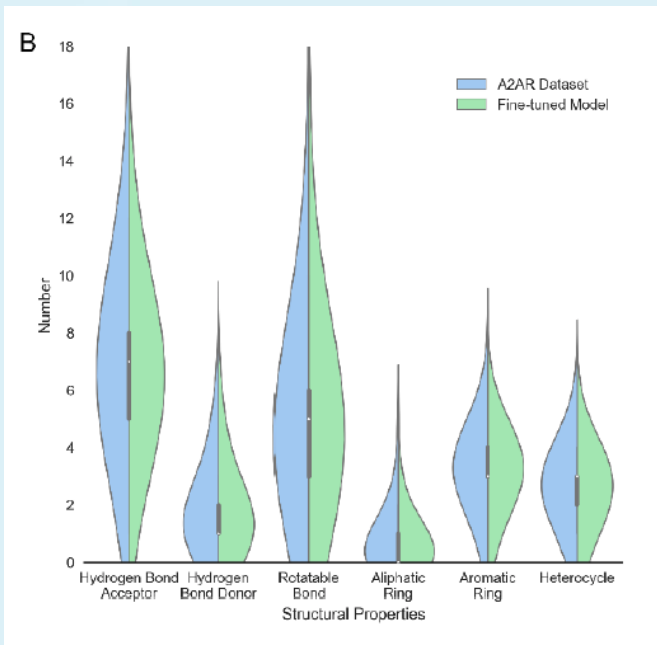


Benchmark on classifying 30 selected bioactive chemicals

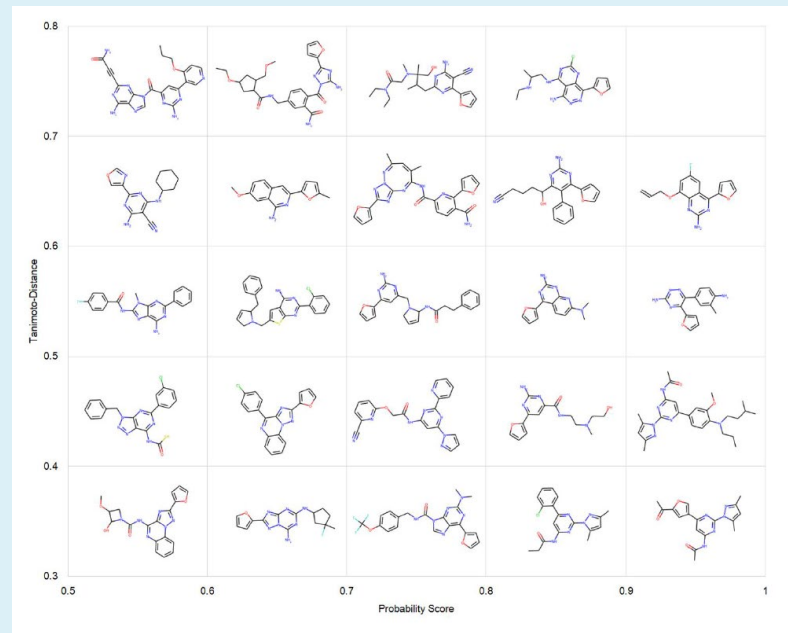
	AUC-ROC	AUC-PR
GAN	<b>0.989±0.007</b>	<b>0.991±0.007</b>
Autoencoder	0.981±0.005	0.985±0.004
cell-profiler	0.967±0.020	0.970±0.019

# An exploration strategy improves the diversity of *de novo* ligands using deep reinforcement learning: a case for the adenosine A<sub>2A</sub> receptor

*De Novo* molecules similar to training molecules

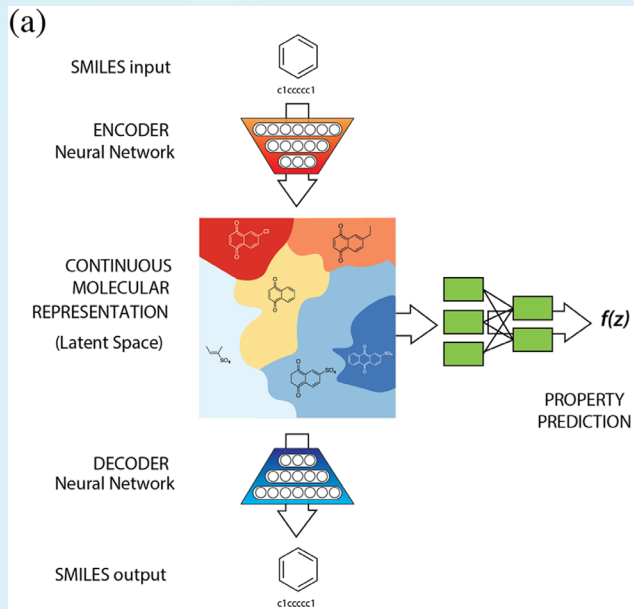


Example generated molecules

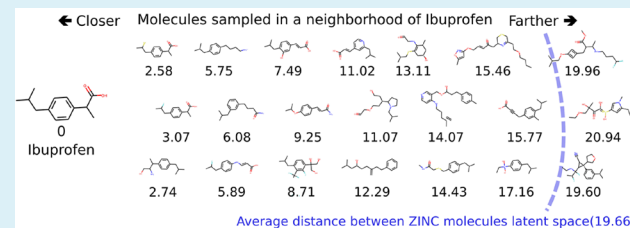


# Automatic Chemical Design Using a Data-Driven Continuous Representation of Molecules

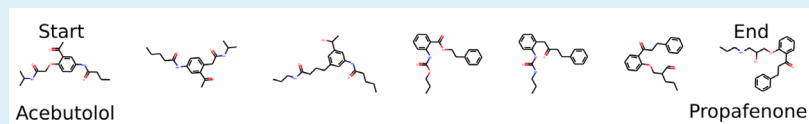
Estimate a continuous latent space for molecular structures by the variational autoencoder



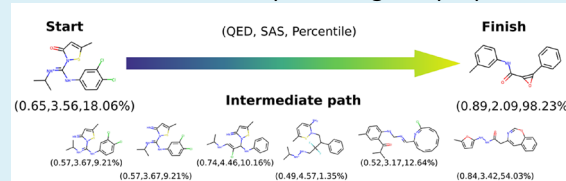
Search similar molecular structures in the latent space



Interpolate between two structures in the latent space

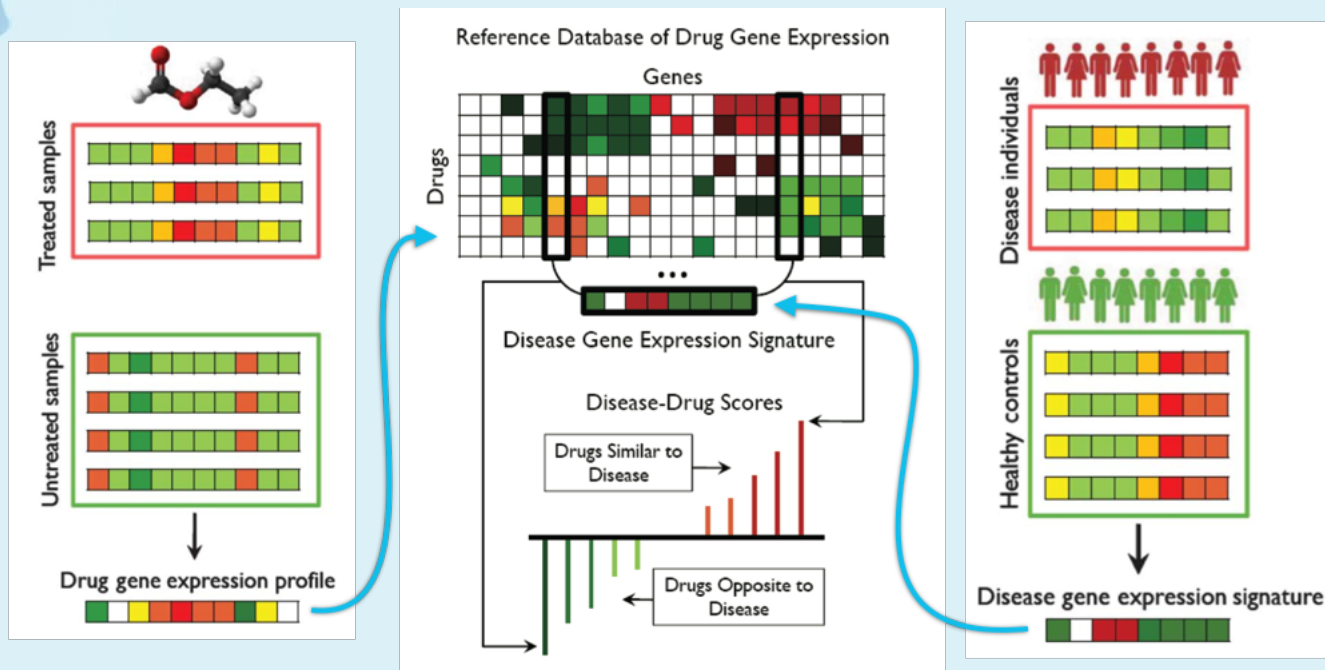


Generate a molecule that optimizes given properties



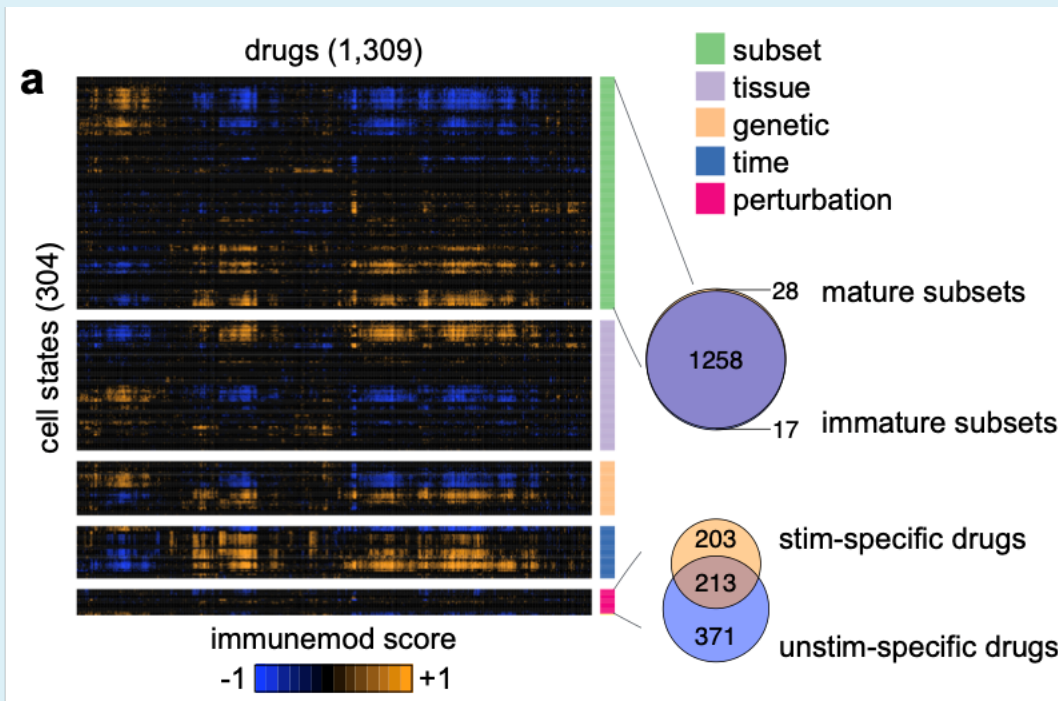
QED: Quantitative Estimation of Drug-likeness  
SAS: Synthetic Accessibility score

# Starting with chemogenomic drug-disease relationships



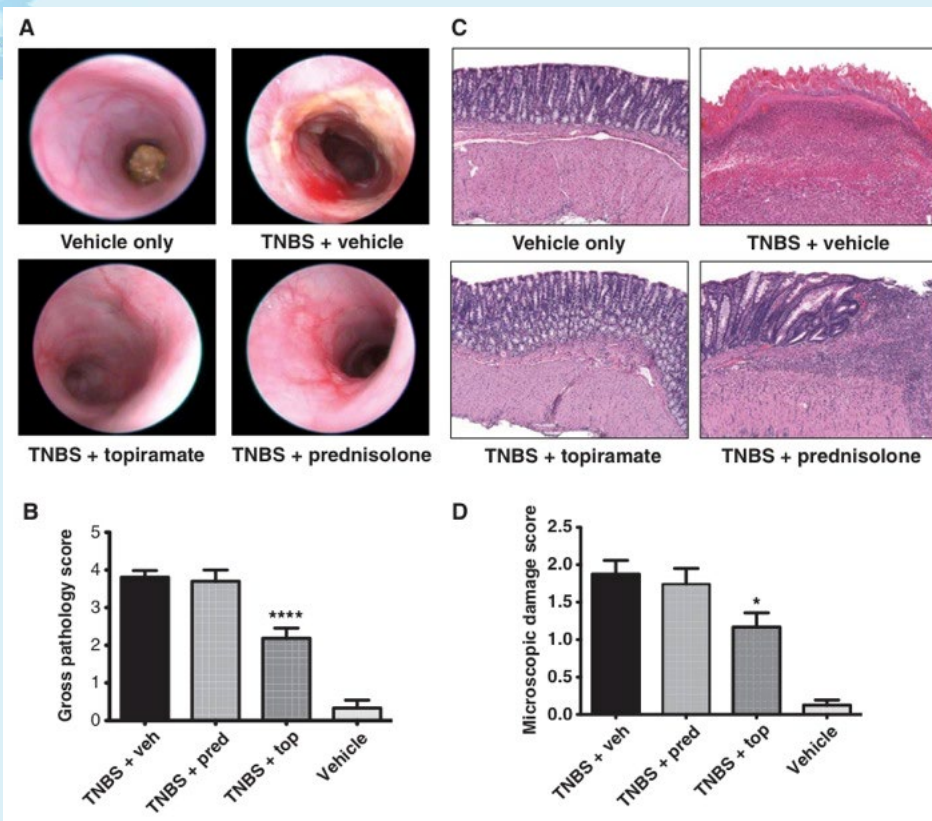
Sirota, M., Dudley, J. T., et al. (2011). Discovery and Preclinical Validation of Drug Indications Using Compendia of Public Gene Expression Data. *Science Translational Medicine*, 3(96).

# Starting with chemogenomic drug-disease relationships





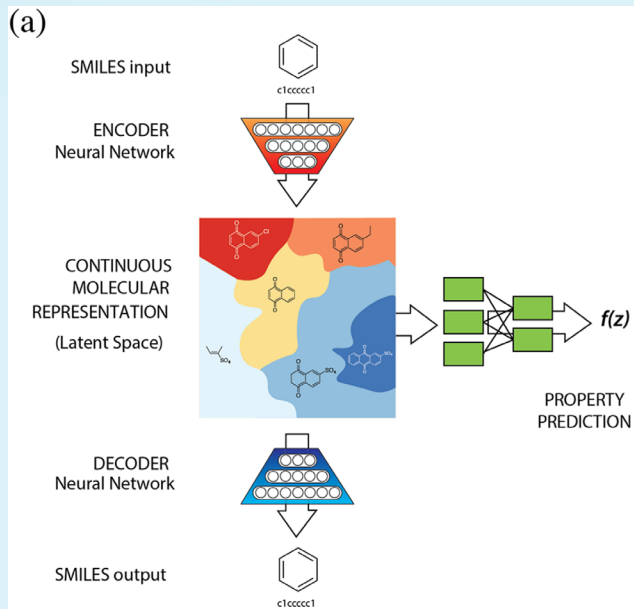
# Starting with chemogenomic drug-disease relationships



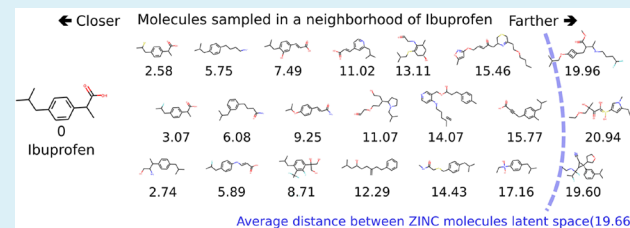
Dudley, J. T., Sirota, M., et al. (2011). Computational Repositioning of the Anticonvulsant Topiramate for Inflammatory Bowel Disease. *Science Translational Medicine*, 3(96).

# Automatic Chemical Design Using a Data-Driven Continuous Representation of Molecules

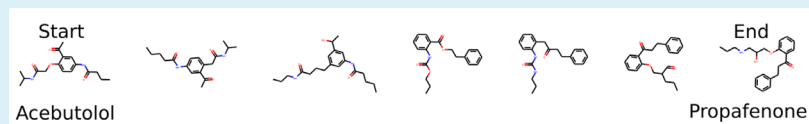
Estimate a continuous latent space for molecular structures by the variational autoencoder



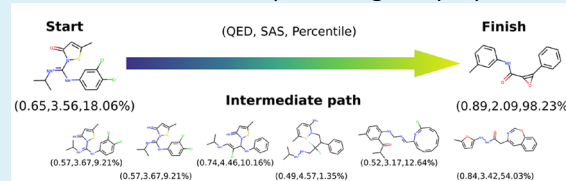
Search similar molecular structures in the latent space



Interpolate between two structures in the latent space

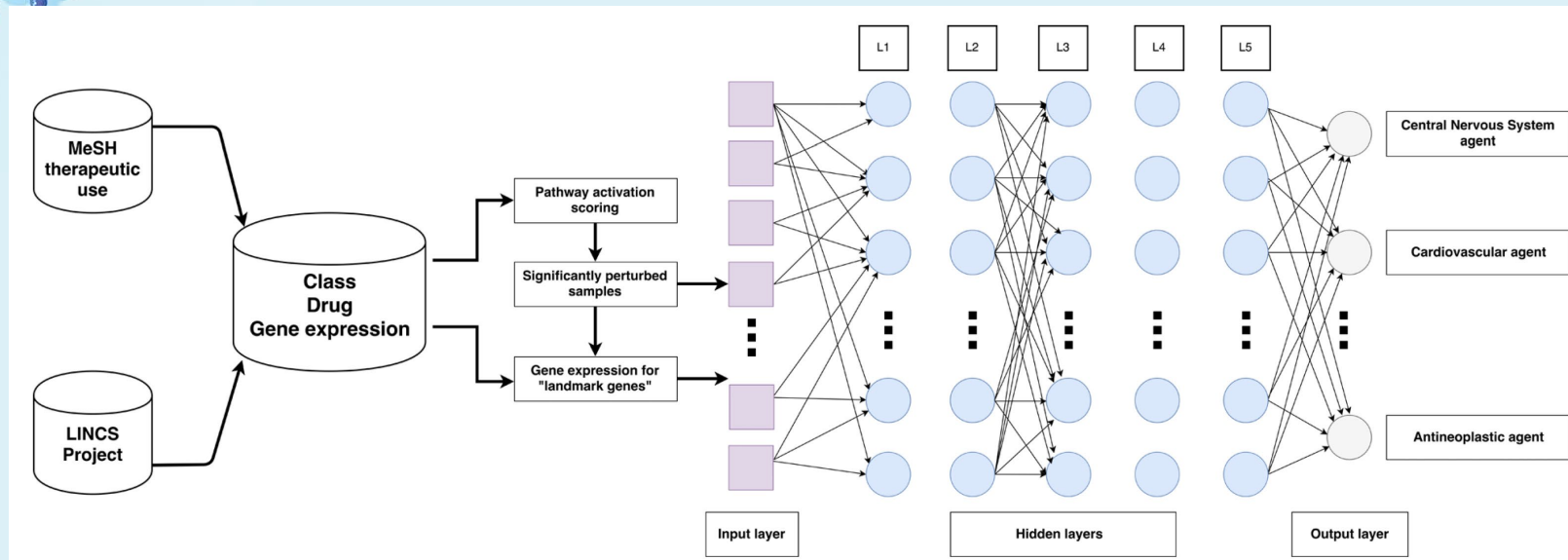


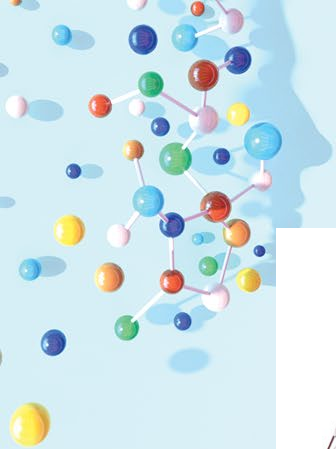
Generate a molecule that optimizes given properties



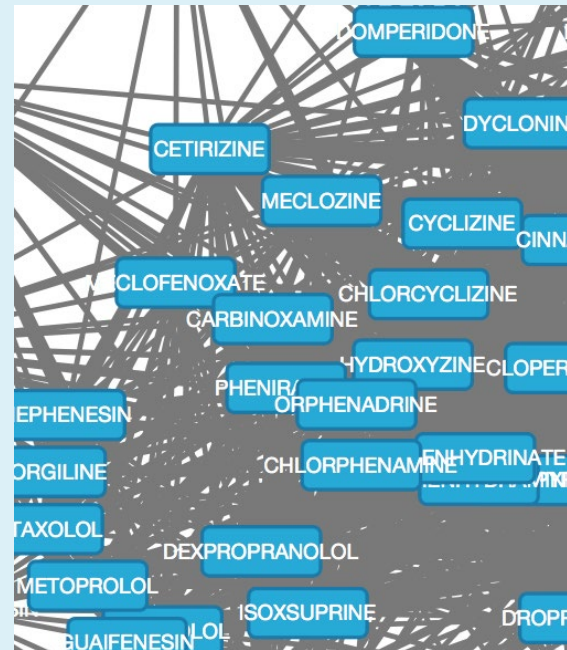
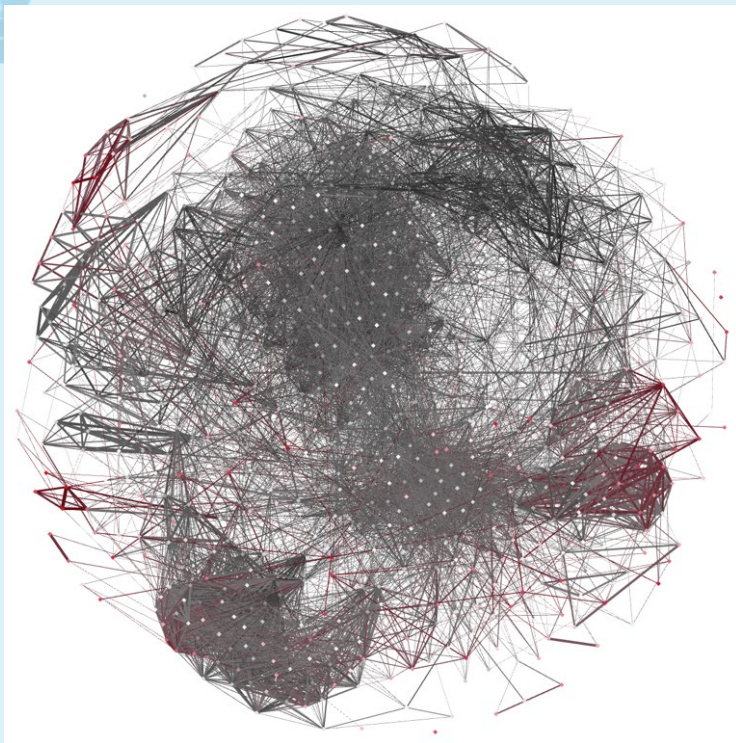
QED: Quantitative Estimation of Drug-likeness  
SAS: Synthetic Accessibility score

# Multi-modal representation learning on chemogenomic data





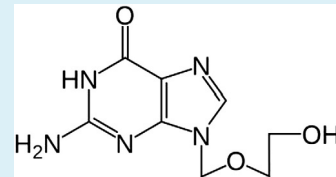
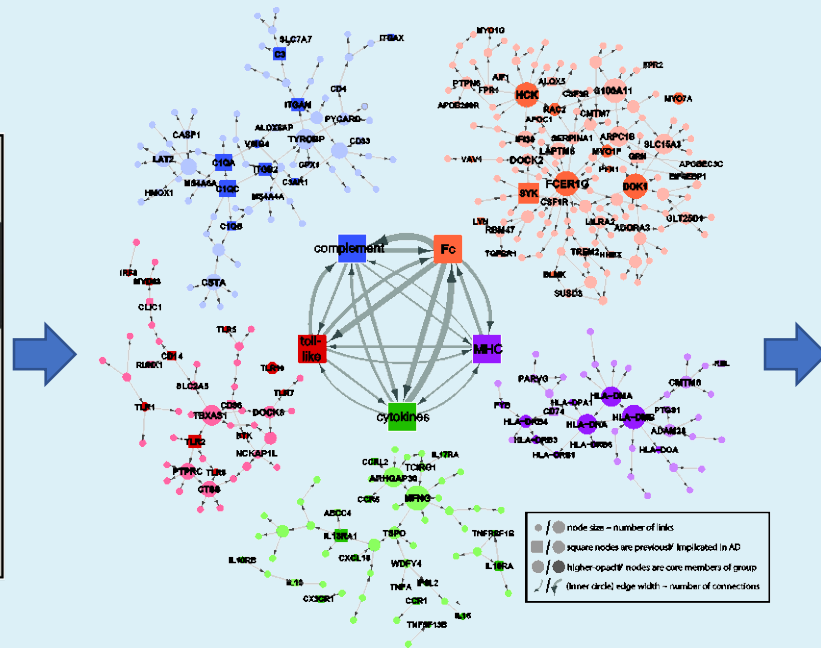
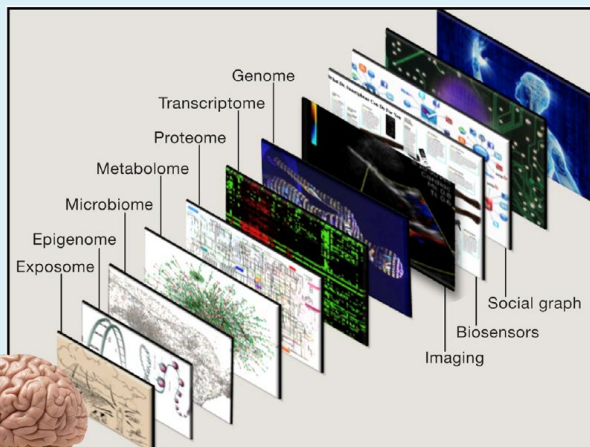
# Multi-modal representation learning on chemogenomic data

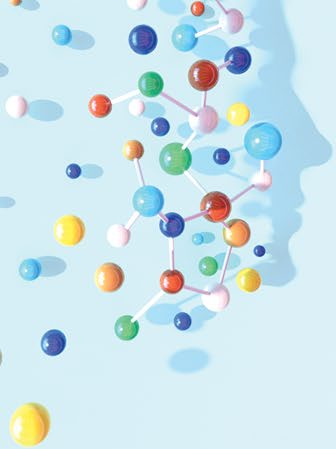


Unpublished work

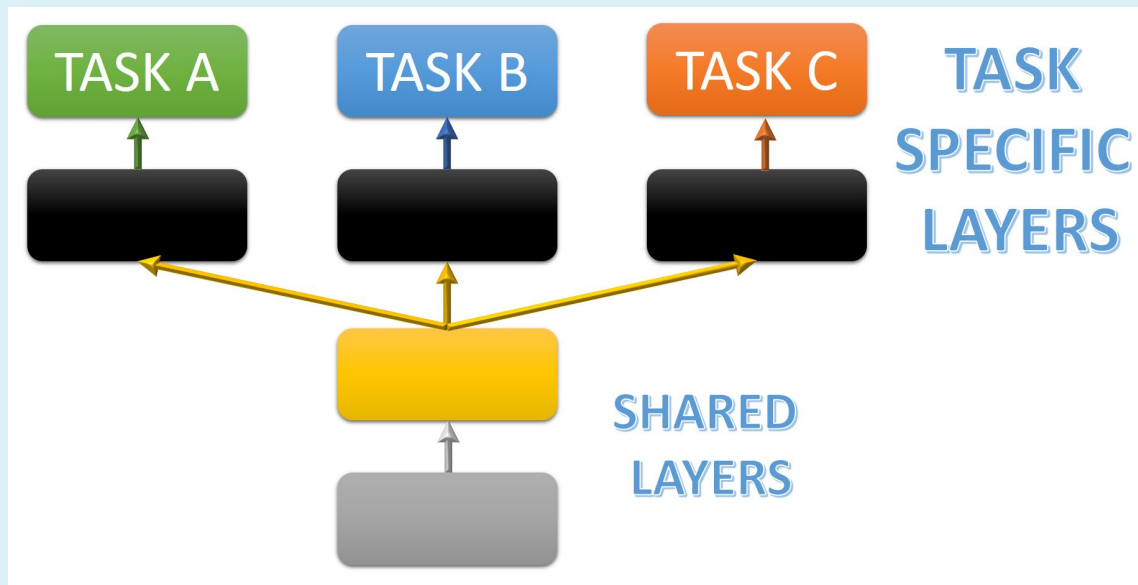


# The holy grail of generative deep learning for drug discovery



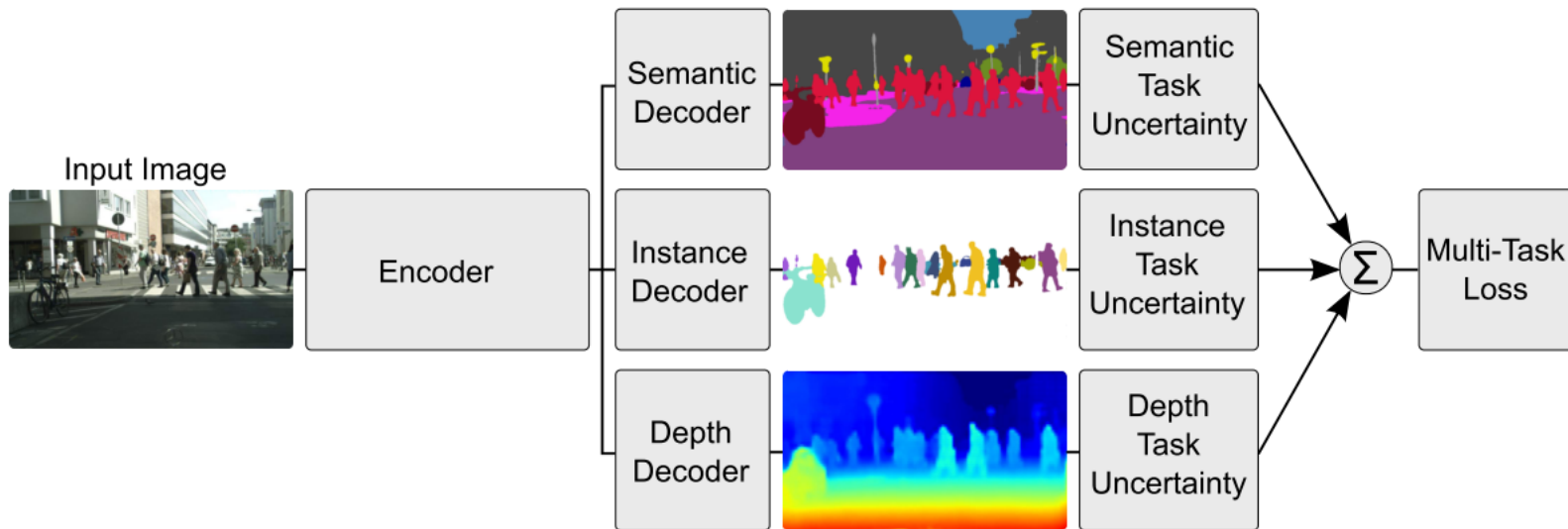


# Multi-Task Learning





# Multi-Task Learning



Source: <http://ruder.io/multi-task/>

# Multi-Task Learning

JOURNAL OF

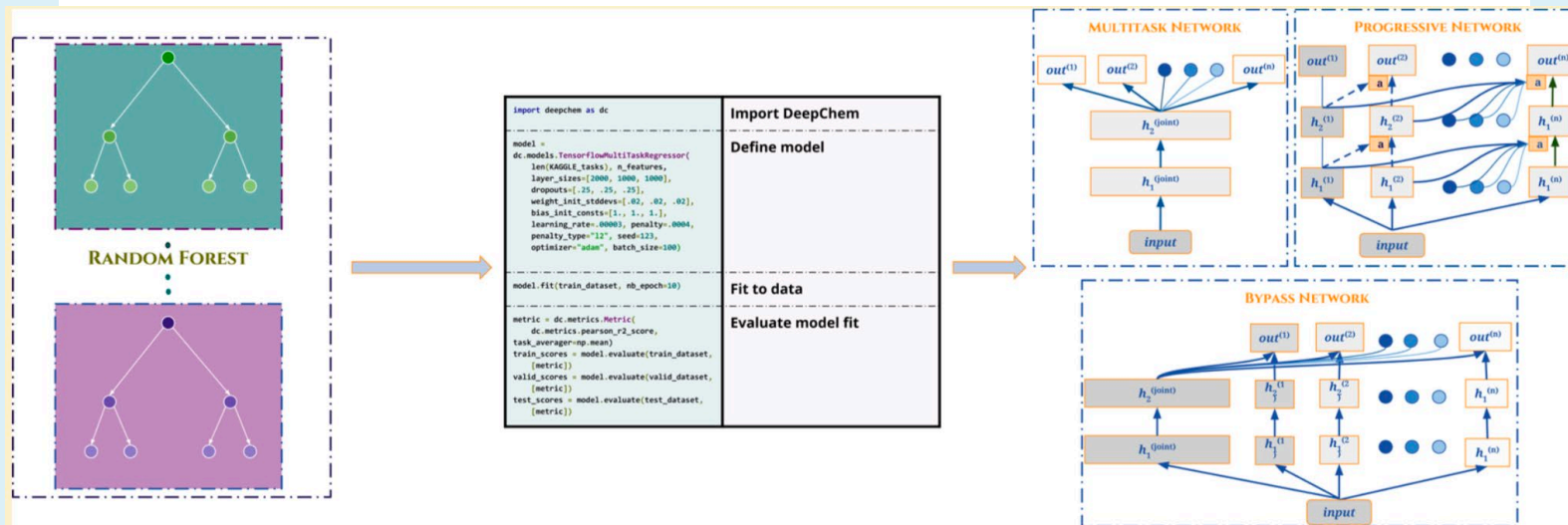
CHEMICAL INFORMATION  
AND MODELING

Article

pubs.acs.org/jcim

## Is Multitask Deep Learning Practical for Pharma?

Bharath Ramsundar,<sup>†</sup> Bowen Liu,<sup>‡</sup> Zhenqin Wu,<sup>‡</sup> Andreas Verras,<sup>¶</sup> Matthew Tudor,<sup>§</sup>  
Robert P. Sheridan,<sup>¶</sup> and Vijay Pande\*,<sup>‡</sup>





# Thank You!

FROM  
MOLECULE TO  
PATIENT

- Mount Sinai
  - Sam Gandy
  - Hao Chi
  - Riccardo Miotto
  - Kipp Johnson
  - Jessica DeFreitas
- ASU
  - Ben Readhead

Email: [joel.dudley@mssm.edu](mailto:joel.dudley@mssm.edu)  
Twitter: [@jdudley](https://twitter.com/jdudley)

