

The rise of deep learning in drug discovery

Russ B Altman, MD, PhD
Stanford University

Machine Learning

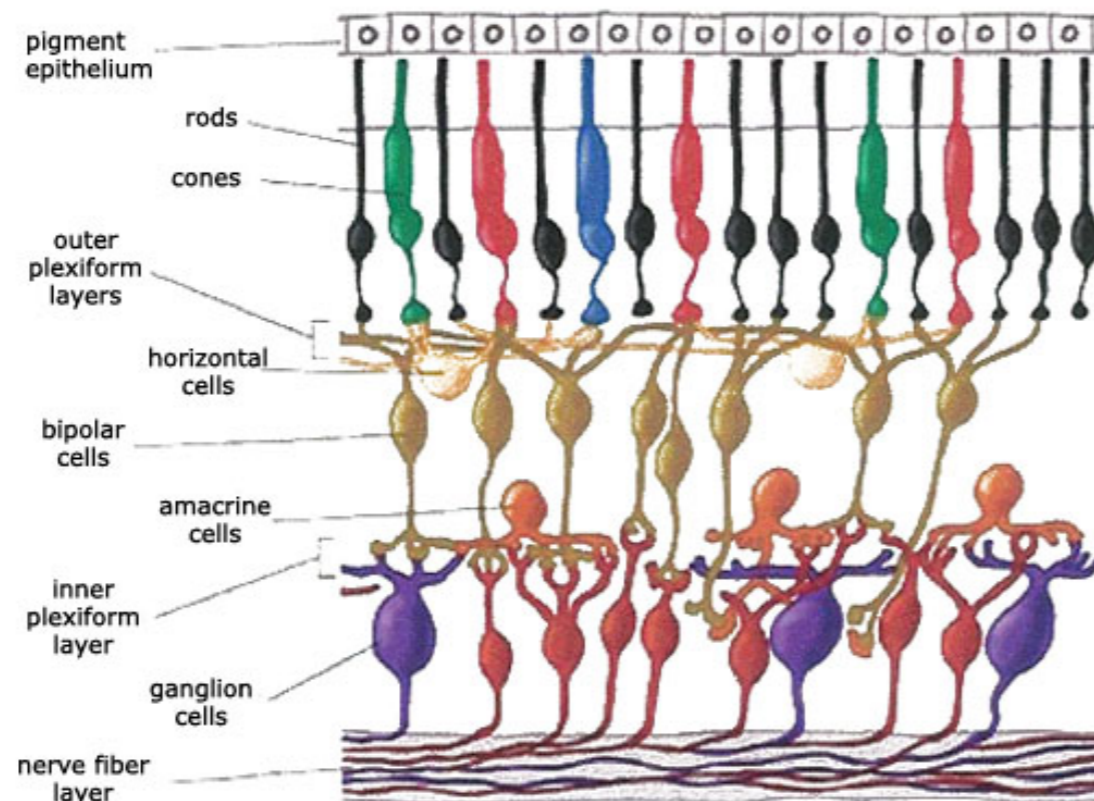
- Machine learning is the discipline within computer science where algorithms are given data and learn how to categorize or classify it.

Two general types of machine learning:

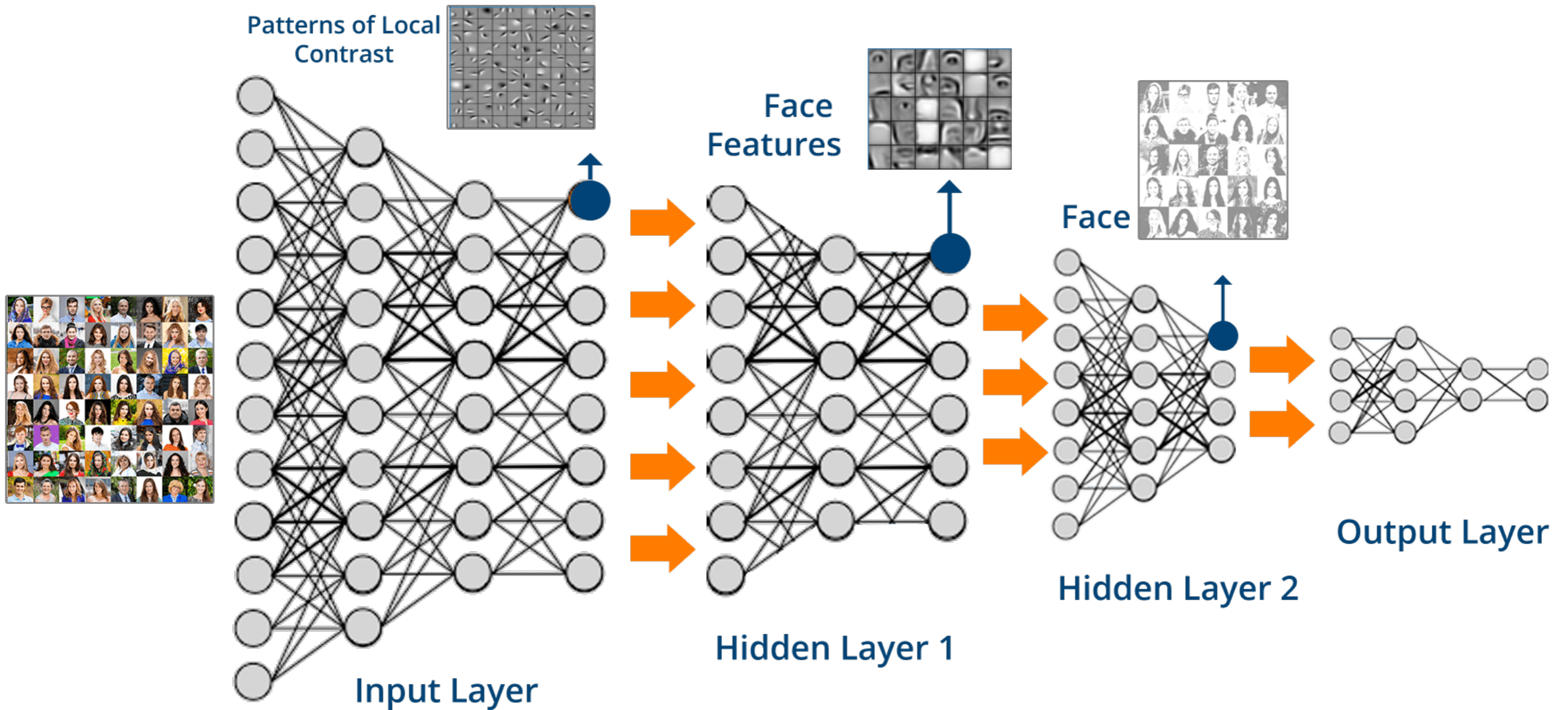
- **Supervised Machine Learning:** provides data as well as labels (the supervision), e.g. lots of data about drugs with labels of whether they are toxic or not
- **Unsupervised Machine Learning:** provides data but no labels, so groups similar objects together = clustering, e.g. lots of data about drugs that can be grouped together based on these data.

Deep Learning

- Deep Learning is a sub-discipline within Machine Learning in which both supervised and unsupervised ML takes place using an analogy to neural processing = neural networks
- Cf. processing of light in the retina.



Deep Learning



Key features of Deep Learning

- Requires typically very large data sets (100K to millions of examples to learn).
 - Some emerging methods (e.g. transfer learning) to reduce the requirement for data
- (Surprisingly) Can learn the key features in the data and does not need human expert to do “feature engineering”
- Given sufficient data, able to learn complex non-linear functions of the input data to predict the output “label”
- If the data is available, the resulting classifiers have outpaced many traditional machine learning approaches by leaps and bounds.

High-performance medicine: the convergence of human and artificial intelligence

Eric J. Topol 

The use of artificial intelligence, and the deep-learning subtype in particular, has been enabled by the use of labeled big data, along with markedly enhanced computing power and cloud storage, across all sectors. In medicine, this is beginning to have an impact at three levels: for clinicians, predominantly via rapid, accurate image interpretation; for health systems, by improving workflow and the potential for reducing medical errors; and for patients, by enabling them to process their own data to promote health. The current limitations, including bias, privacy and security, and lack of transparency, along with the future directions of these applications will be discussed in this article. Over time, marked improvements in accuracy, productivity, and workflow will likely be actualized, but whether that will be used to improve the patient-doctor relationship or facilitate its erosion remains to be seen.

Medicine is at the crossroad of two major trends. The first is a failed business model, with increasing expenditures and jobs allocated to healthcare, but with deteriorating key outcomes, including reduced life expectancy and high infant, childhood, and maternal mortality in the United States^{1,2}. This exemplifies a paradox that is not at all confined to American medicine: investment of more human capital with worse human health outcomes. The second is the generation of data in massive quantities, from sources such as high-resolution medical imaging, biosensors

type of medical scan, with more than 2 billion performed worldwide per year. In one study, the accuracy of one algorithm, based on a 121-layer convolutional neural network, in detecting pneumonia in over 112,000 labeled frontal chest X-ray images was compared with that of four radiologists, and the conclusion was that the algorithm outperformed the radiologists. However, the algorithm's AUC of 0.76, although somewhat better than that for two previously tested DNN algorithms for chest X-ray interpretation³, is far from optimal. In addition, the test used in this study is not necessarily comparable

PubMed ID: 30617339

Areas of Deep Learning Success

- **Imaging:** (the first application area where it emerged successful) Applications in radiology, pathology, dermatology. Images are a 2D array of data.
- **Text processing:** Application in mining the published literature, mining clinical notes, taking dictation, language translation. Text is a 1D array of data.

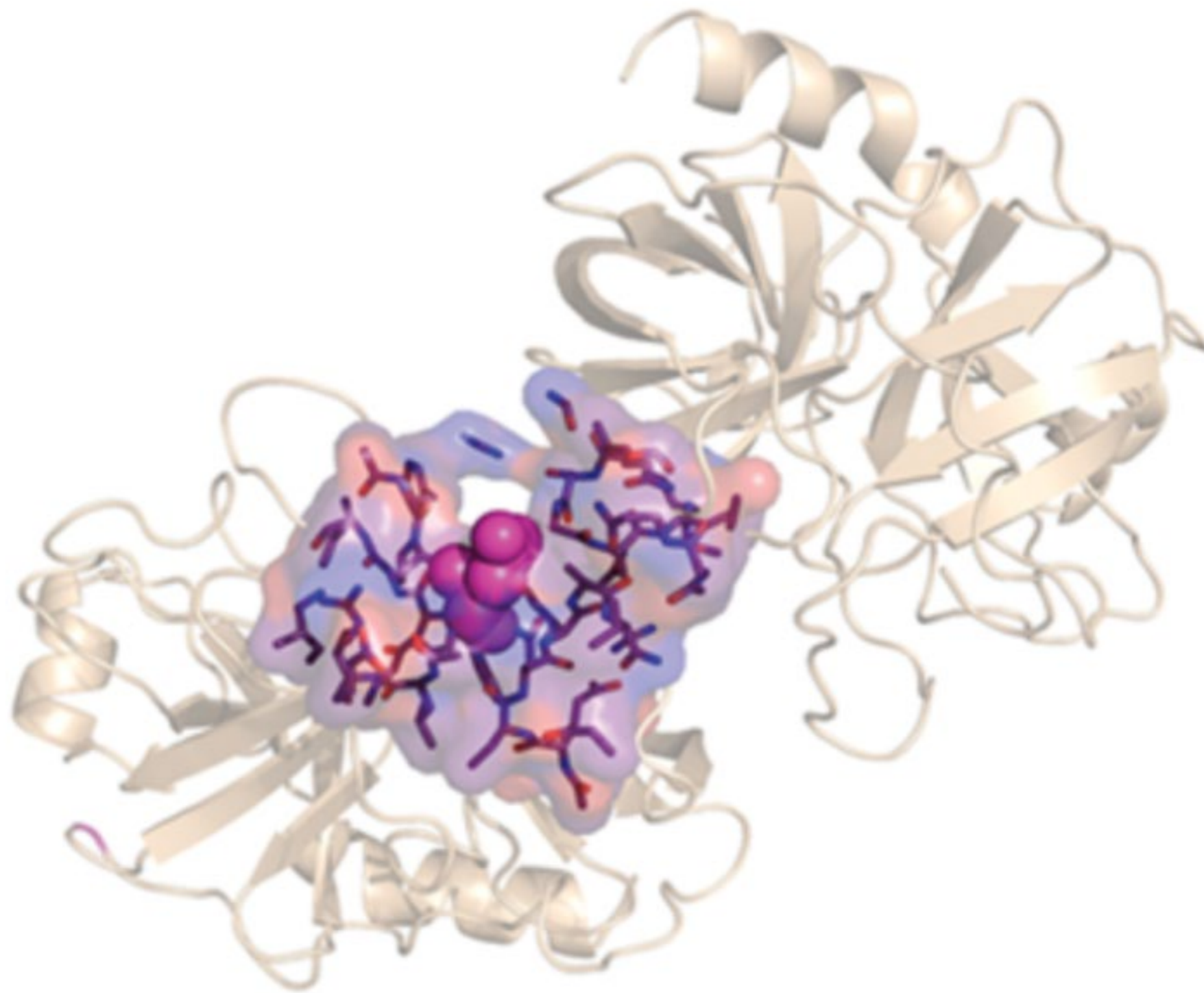
(Progress in DNA/Protein sequence analysis by analogy)

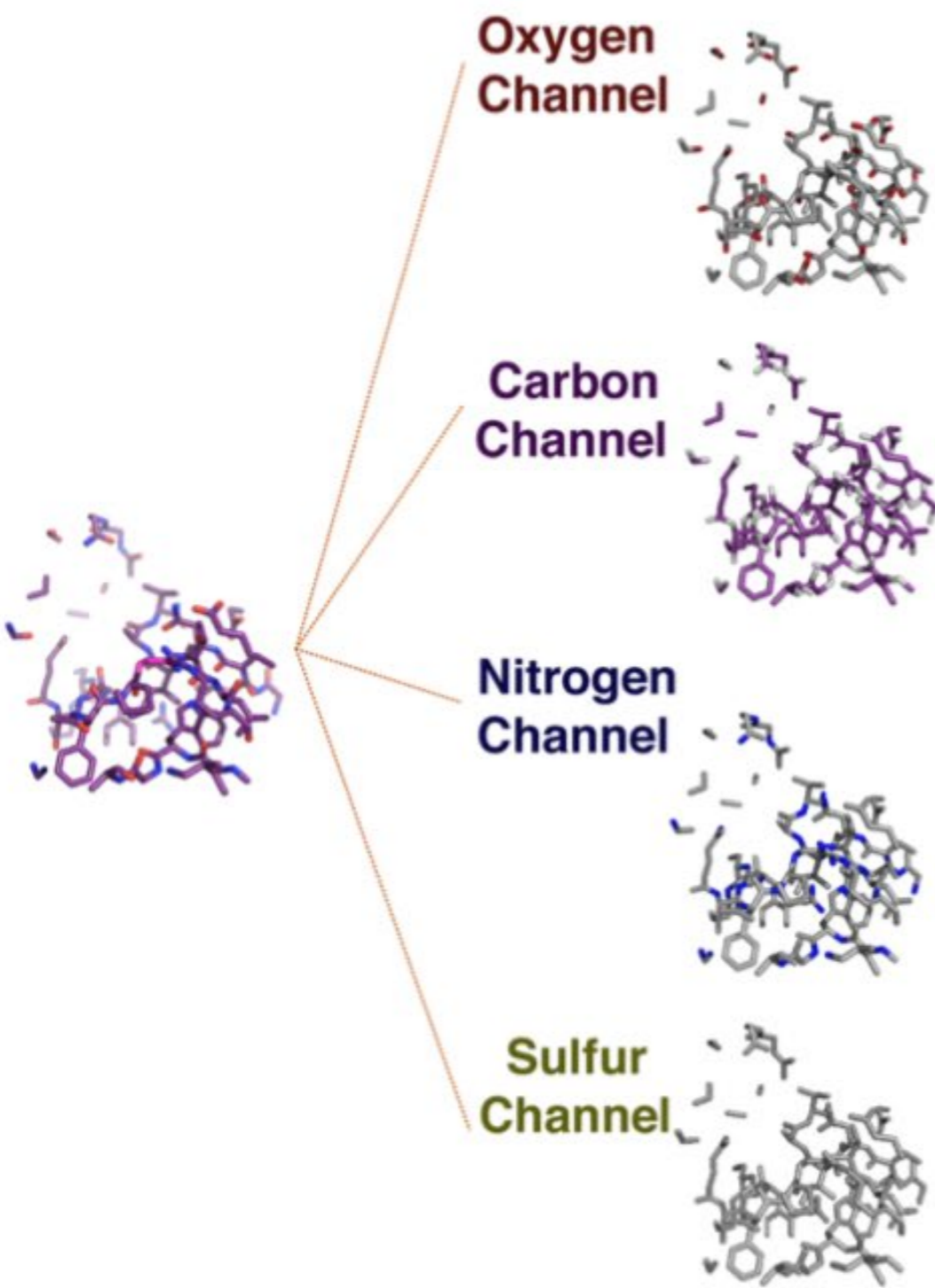
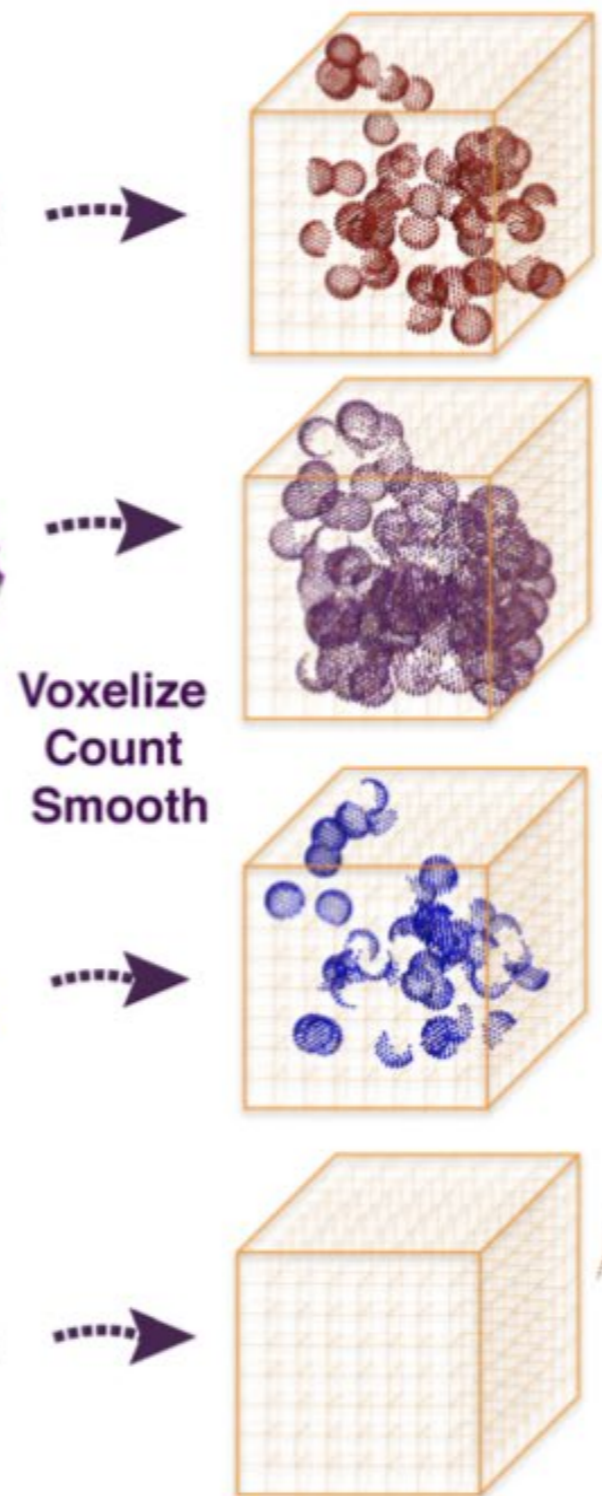
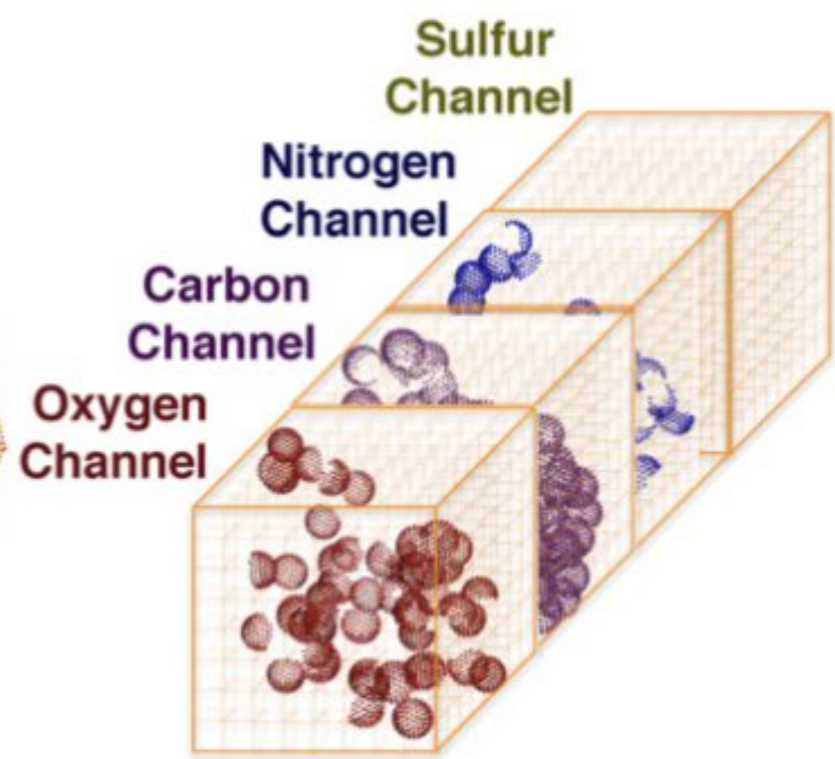
More challenging: integrating data that is not a regular array of homogenous data.

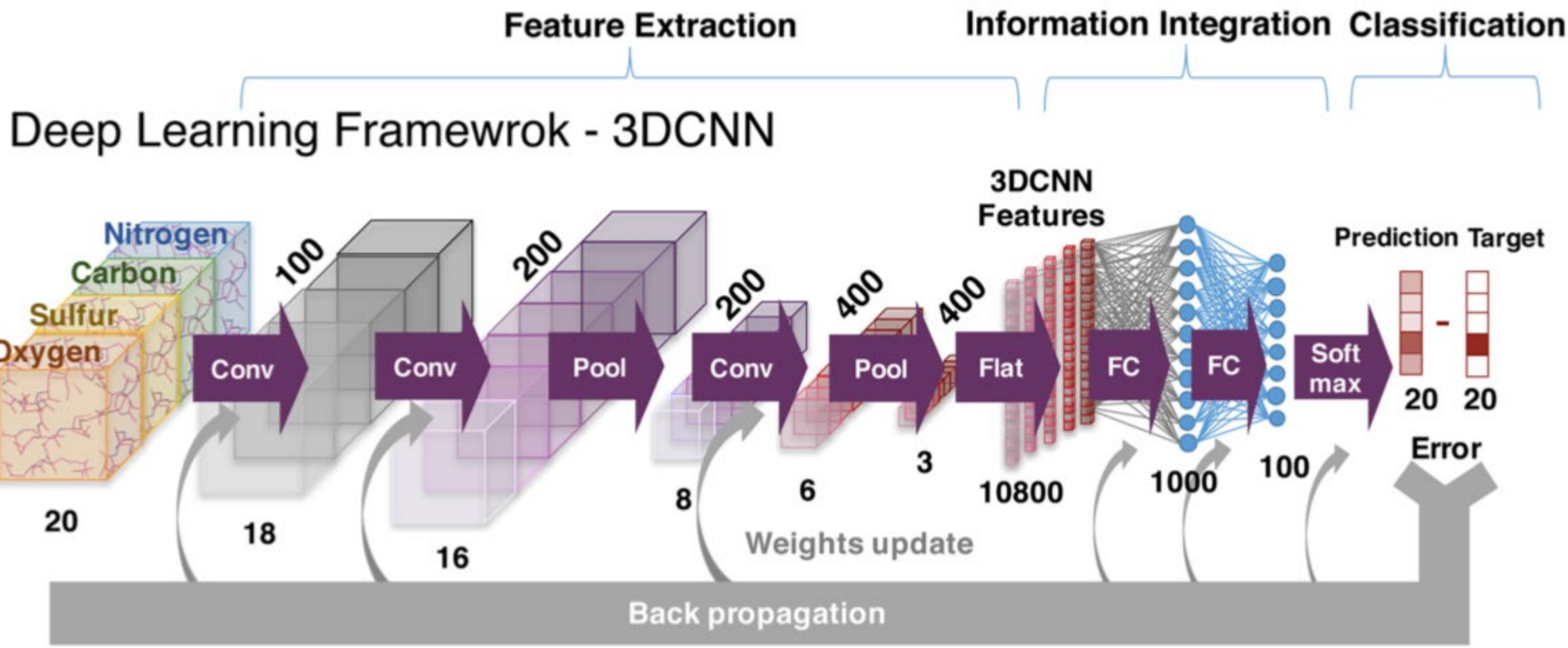
Sample applications in drug discovery

- **3D structure-based discovery/screening of ligands**
- **Cellular networks to find targets**
- **Gene-drug-phenotype networks to predict toxicity**
- **Predicting gene function (phenotype) from genotype**

Deep Networks for Understanding Molecular Structure

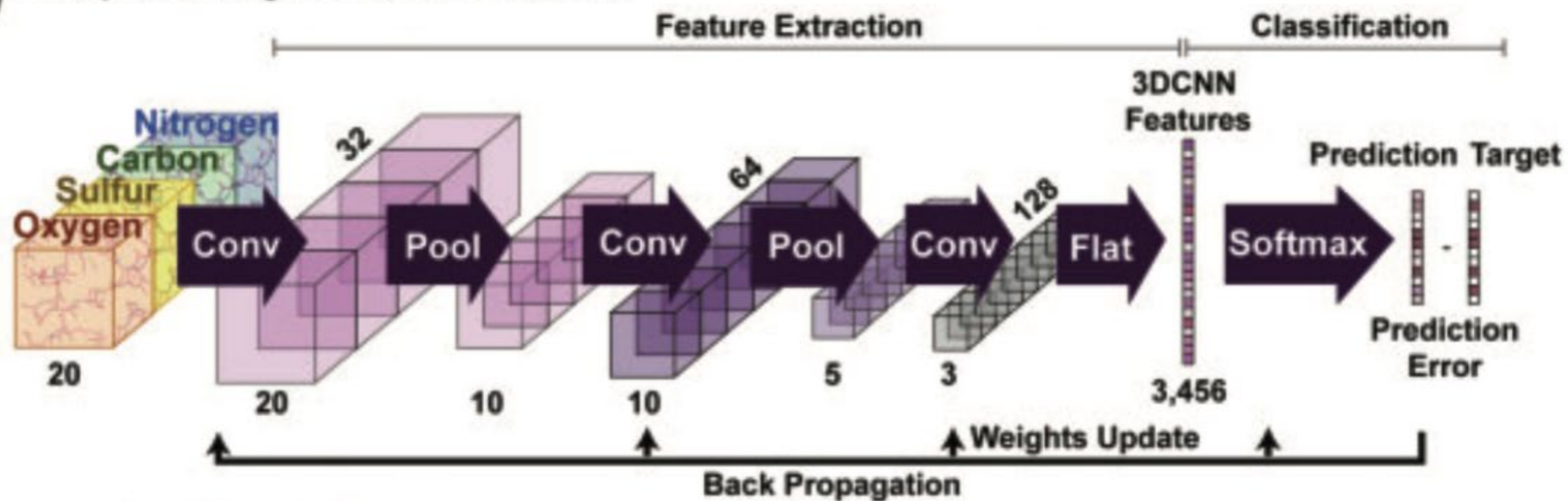


a**b****c**

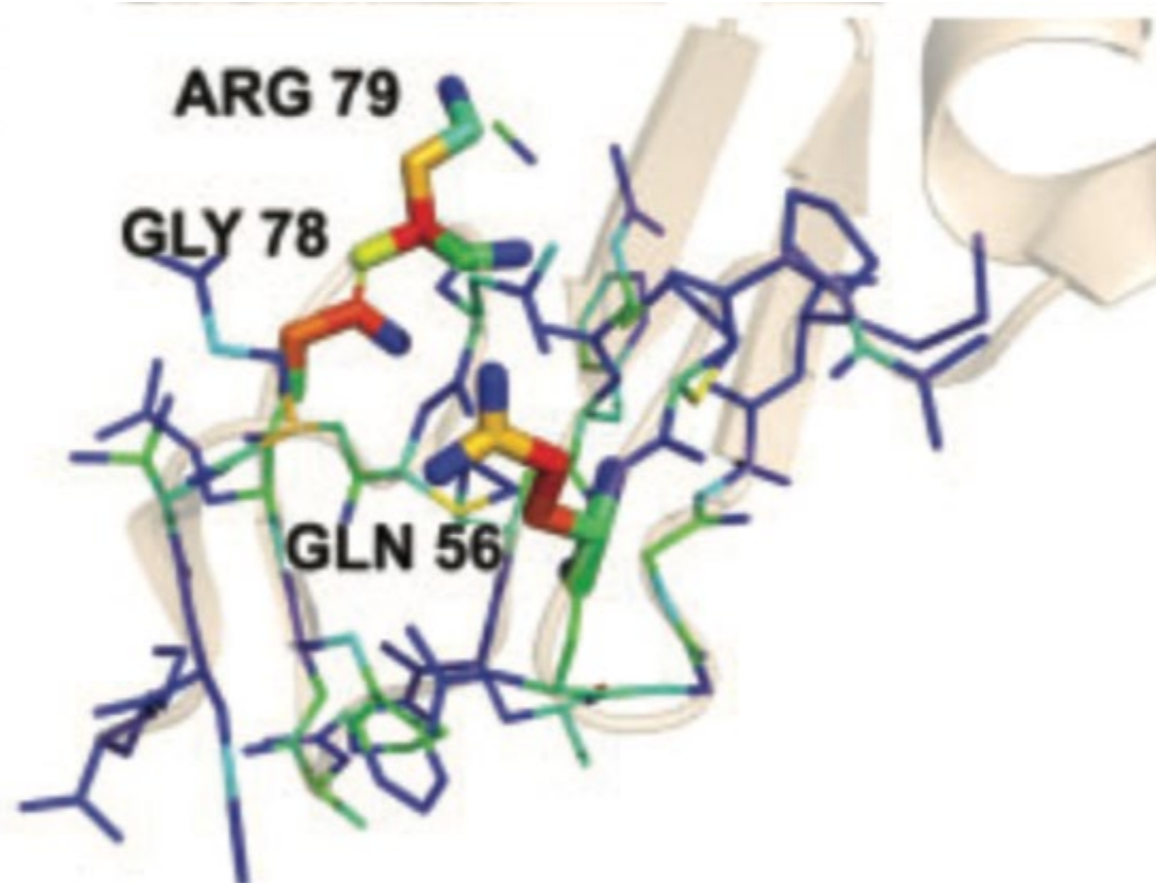


Recognizing Epidermal Growth Factor Binding Site

(a) Deep Learning Framework - 3DCNN



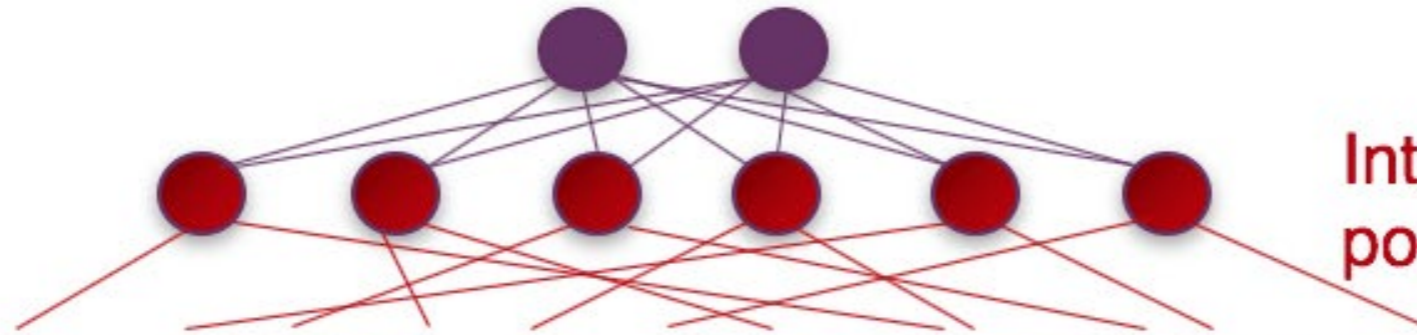
(c)



A Supervised Classifier for Binding Prediction

Binding Prediction

Softmax Classifier



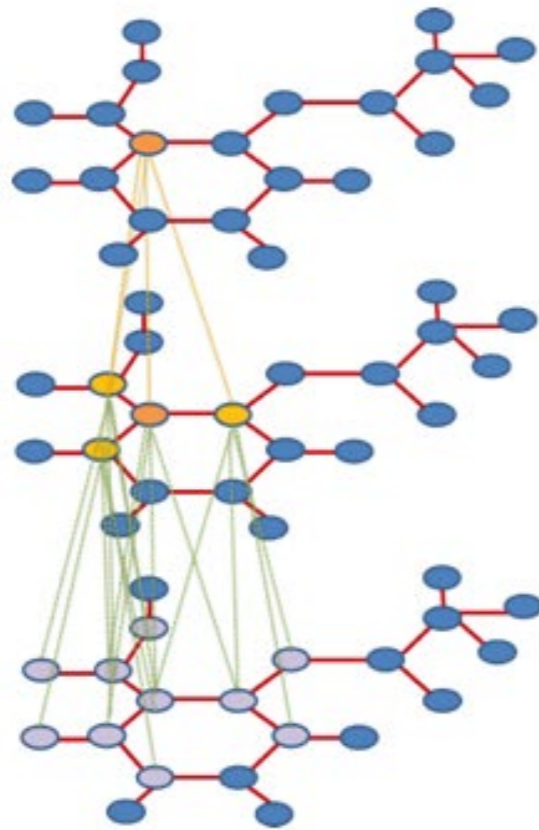
Interactions between pockets and molecules

SMALL MOLECULE REPRESENTATION

PROTEIN POCKET REPRESENTATION

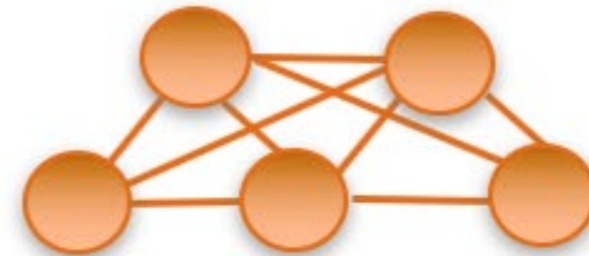


Graph Convolution



Duvenaud, David K., et al. 2015

Graph Convolution



Amino Acid Environments



Deep Networks for predicting new drug targets

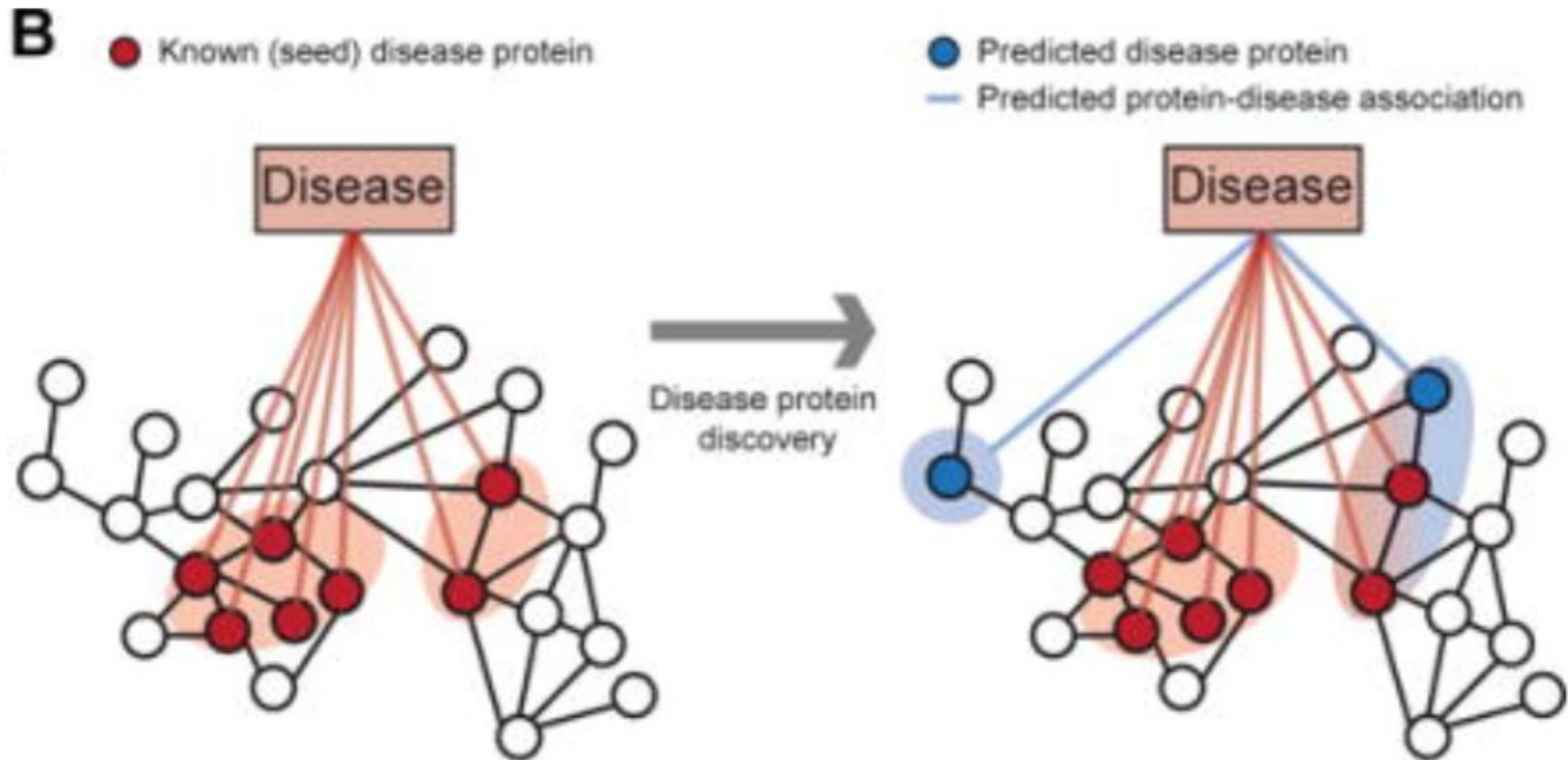
**Same Network
analyzed by two Deep
Learning Methods.**

**Top shows
neighborhoods
based on connectivity
of nodes
in network**



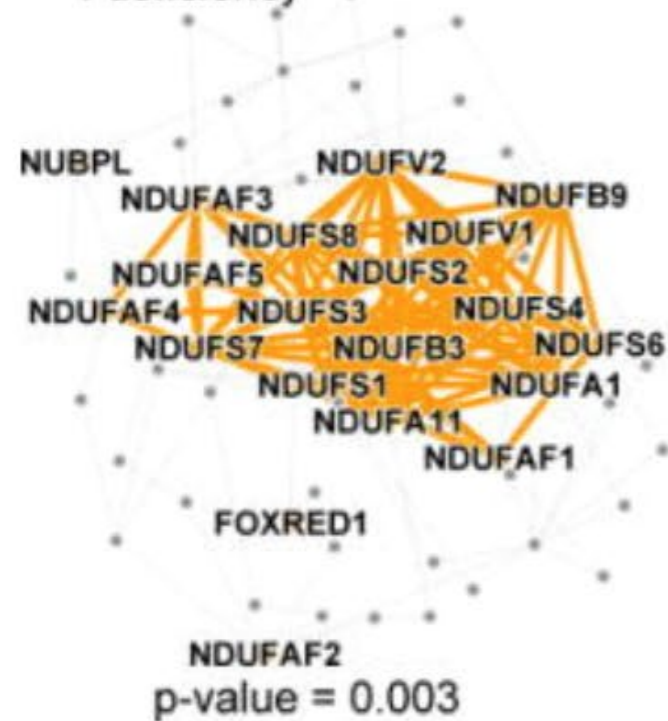
**Bottom shows nodes
that are
similar based on the
role that they
play in the network**



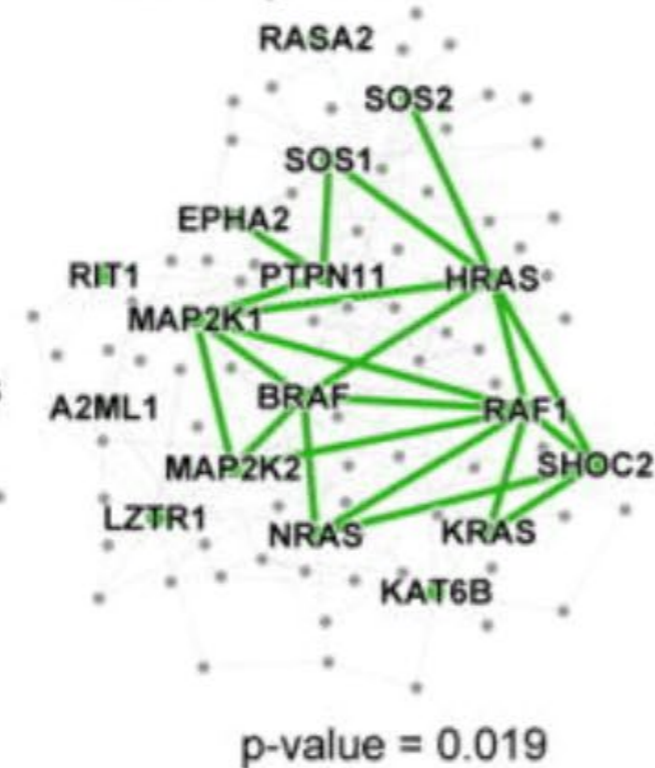


Agrawal M, Zitnik M, Leskovec J. Large-scale analysis of disease pathways in the human interactome. Pac Symp Biocomput. 2018;23:111-122. PMID: 29218874

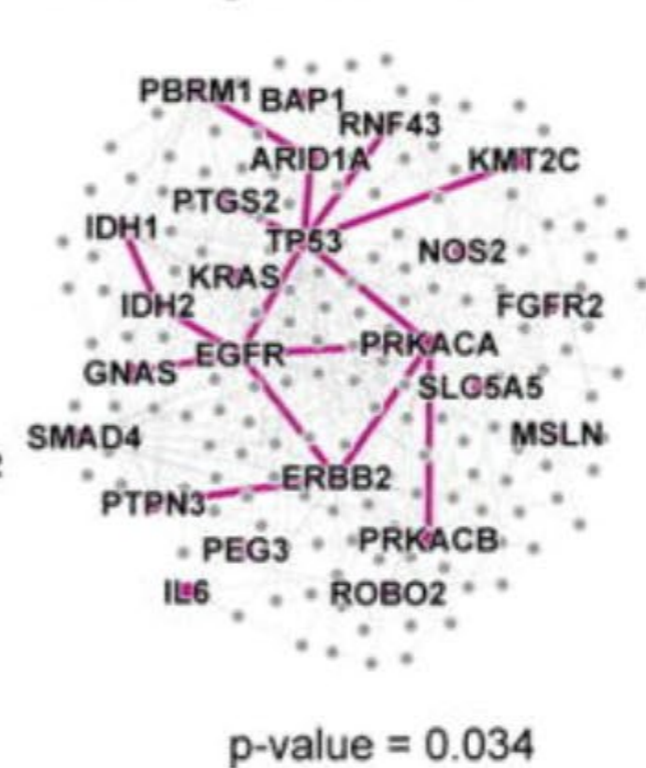
A Mitochondrial complex I deficiency



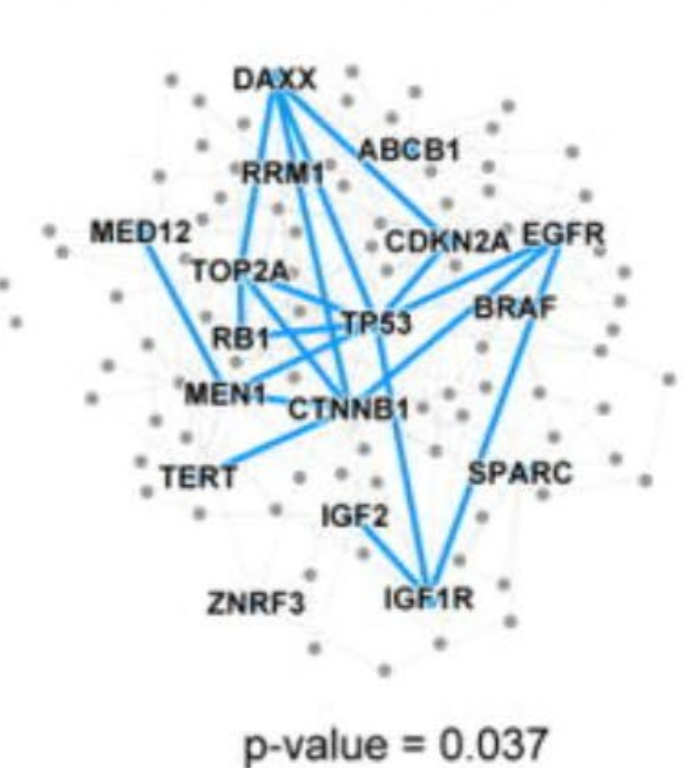
B Noonan syndrome



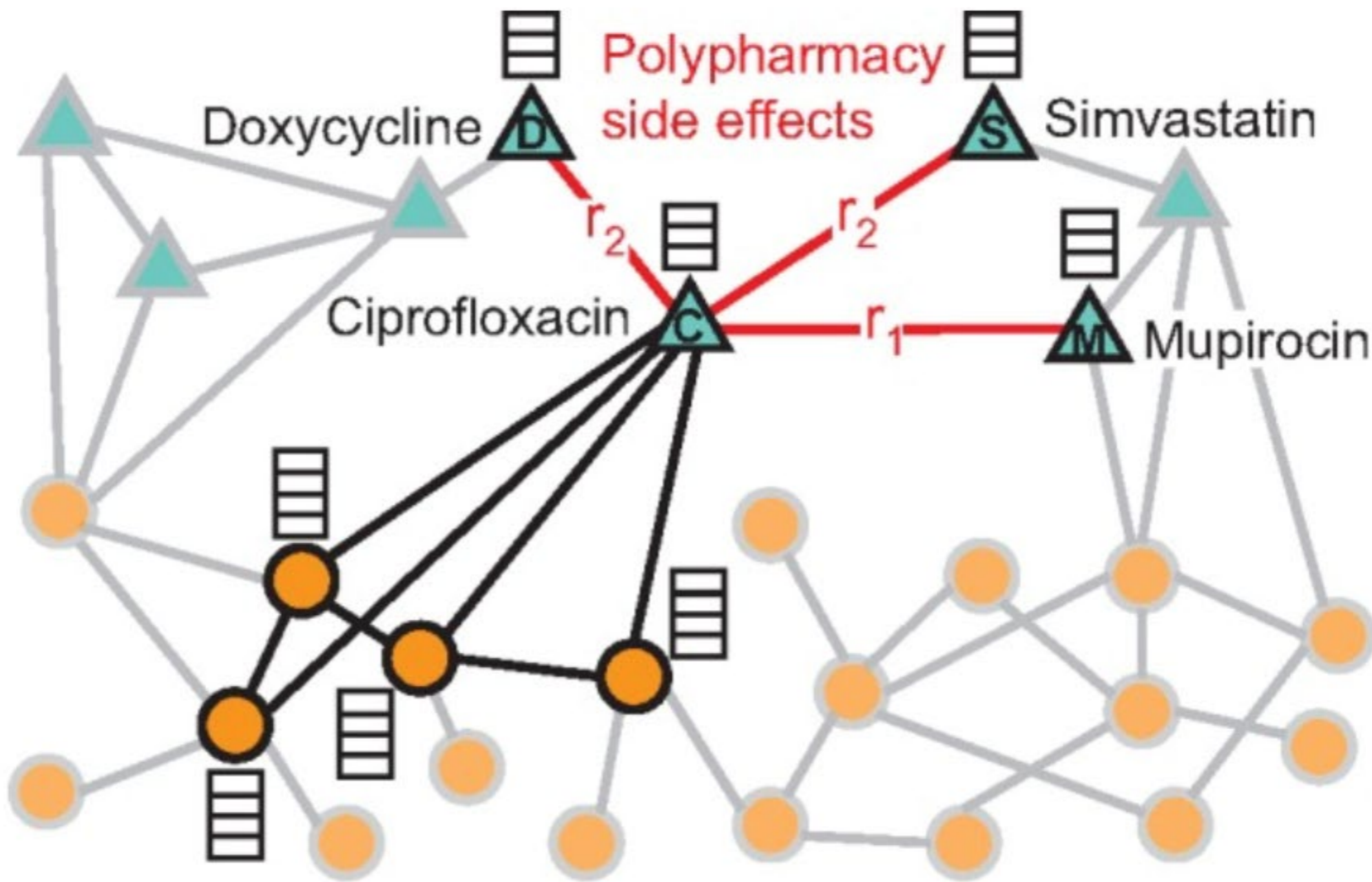
C Cholangiocarcinoma



D Adrenal cortex carcinoma



Deep Networks for predicting drug-drug interactions



▲ Drug ● Protein

r_1 Gastrointestinal bleed side effect

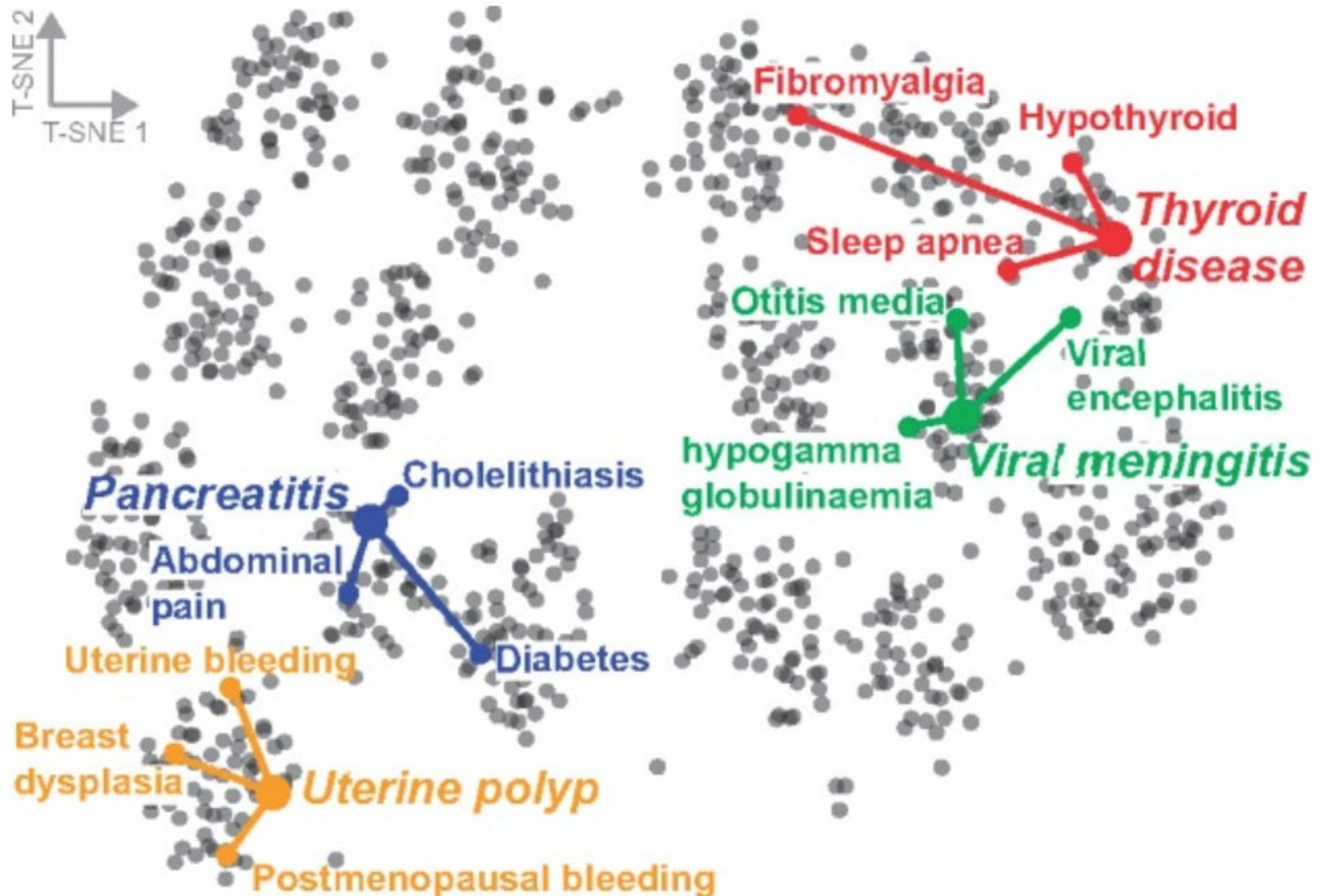
r_2 Bradycardia side effect

Node feature vector

▲—● Drug-protein interaction

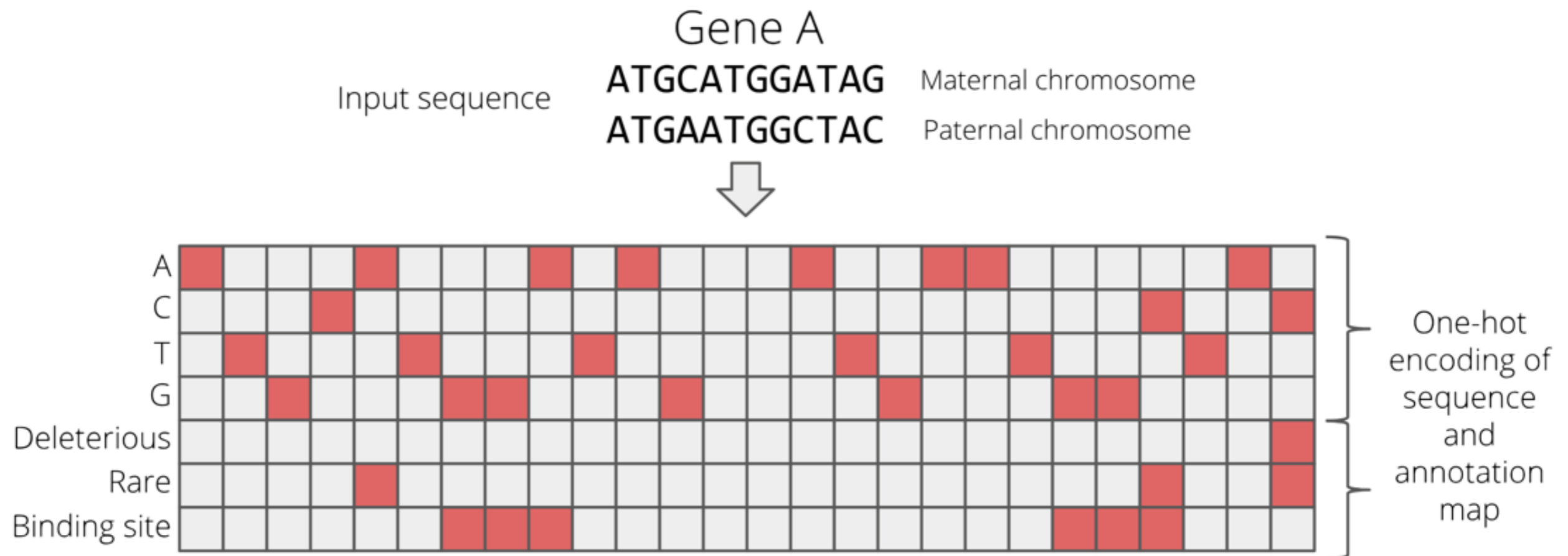
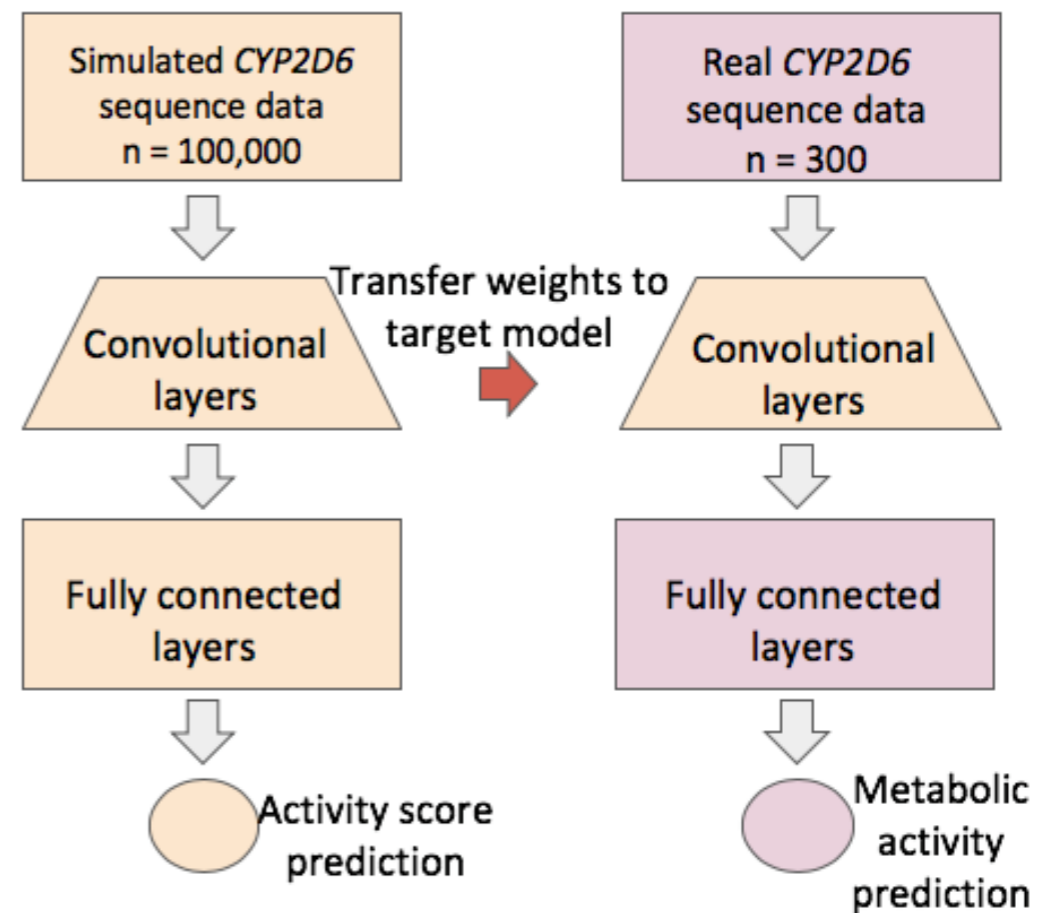
●—● Protein-protein interaction

Co-occurring side effects in drug combinations

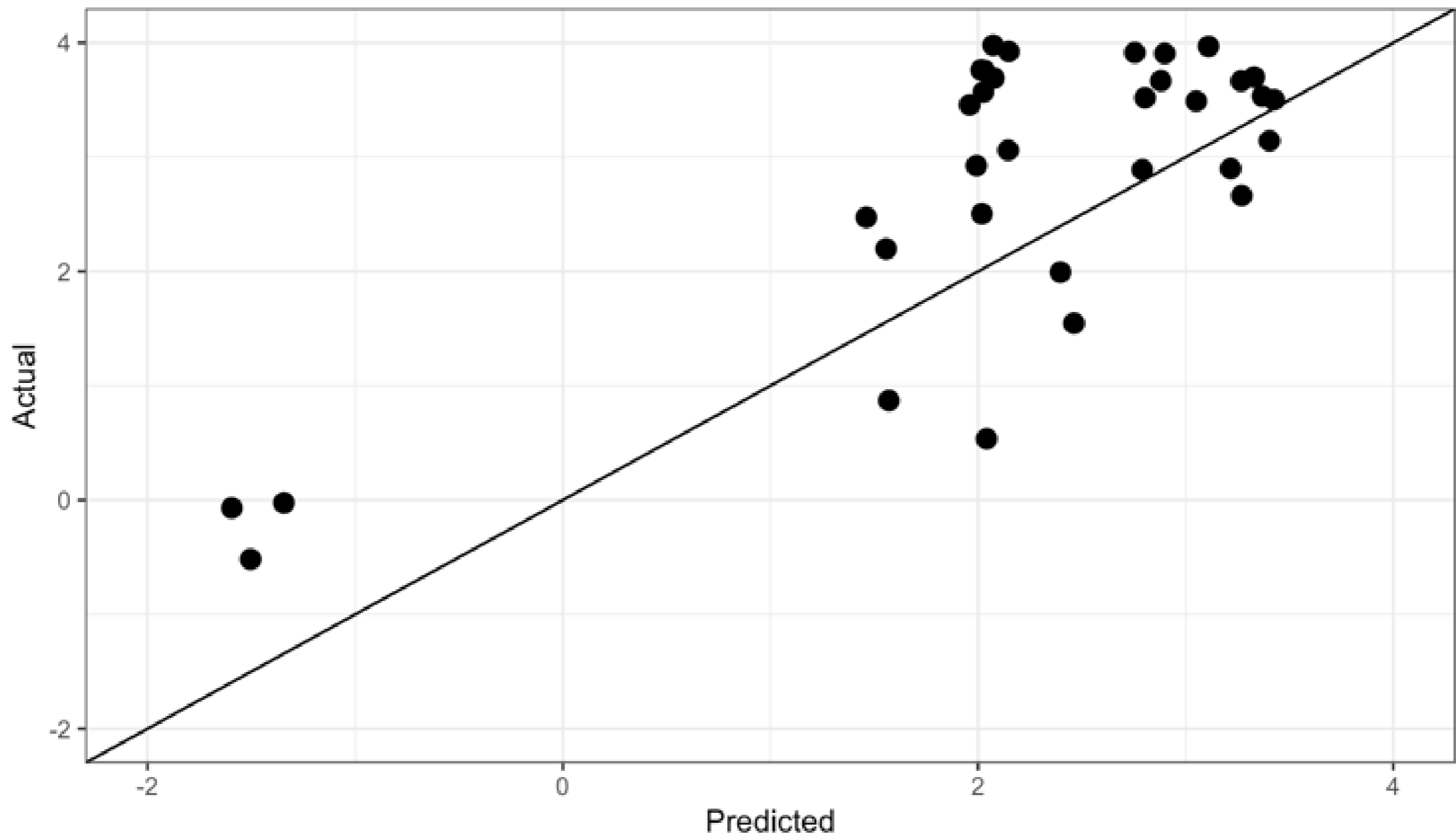


Deep Networks for predicting protein function

General strategy for CYP2D6 prediction and representation of DNA sequence data



Post-transfer Predicted Values



Greg McInnes, unpublished results (with Erika Woodahl = U. Montana)

Summary

- Deep Learning is a type of machine learning
 - Heavy data requirements
 - Able to fit data extremely well
 - Need to validate rigorously with held-out data
- Early applications show promise for data-driven drug discovery, improving speed and quality of nominated drug targets.

Thanks!

russ.altman@stanford.edu

Support: NIH NIGMS, NIH NLM, NIH NCATS,
FDA CERSI, Chan-Zuckerberg Biohub

