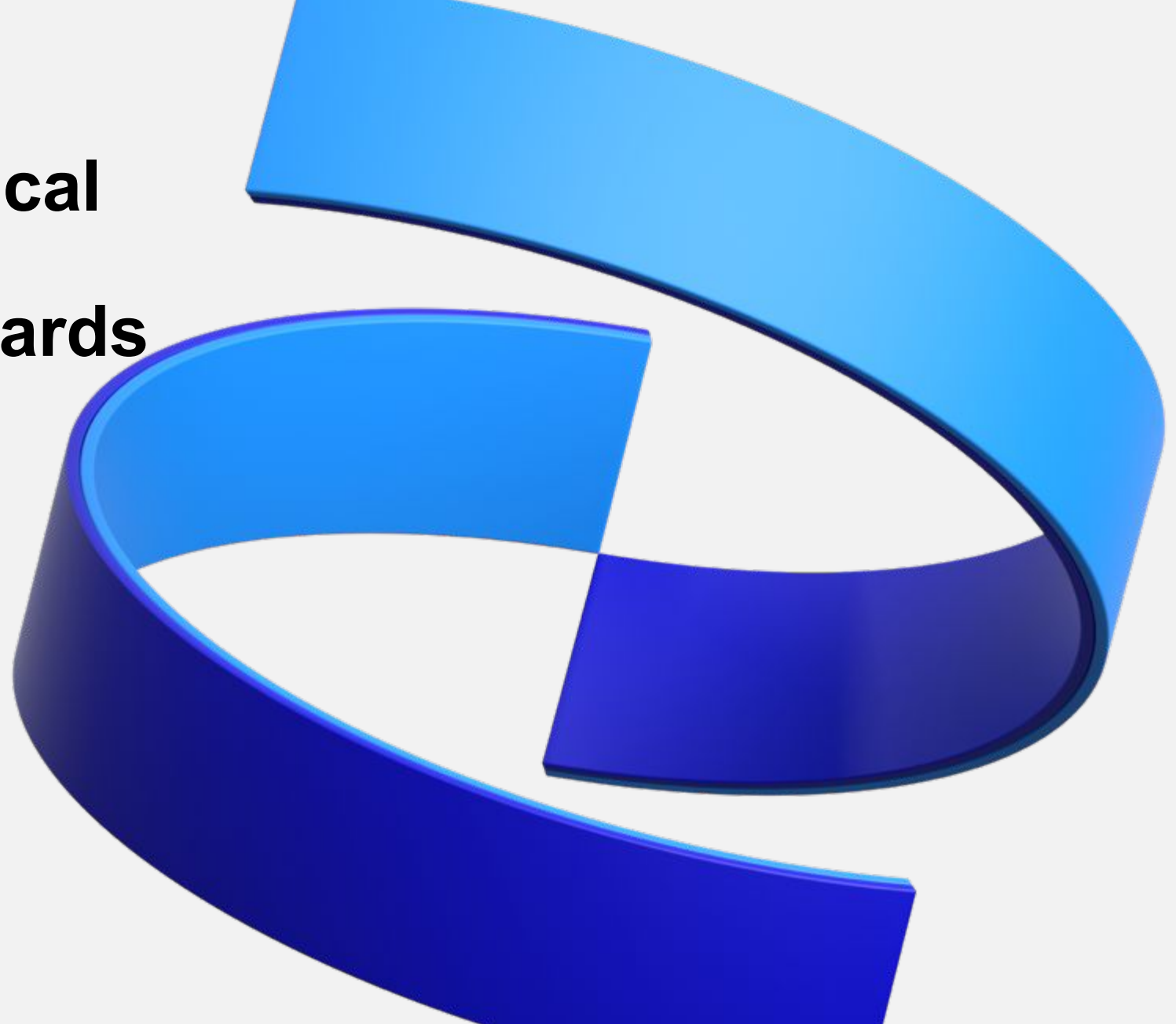# Applications of classical and contemporary machine learning towards drug discovery
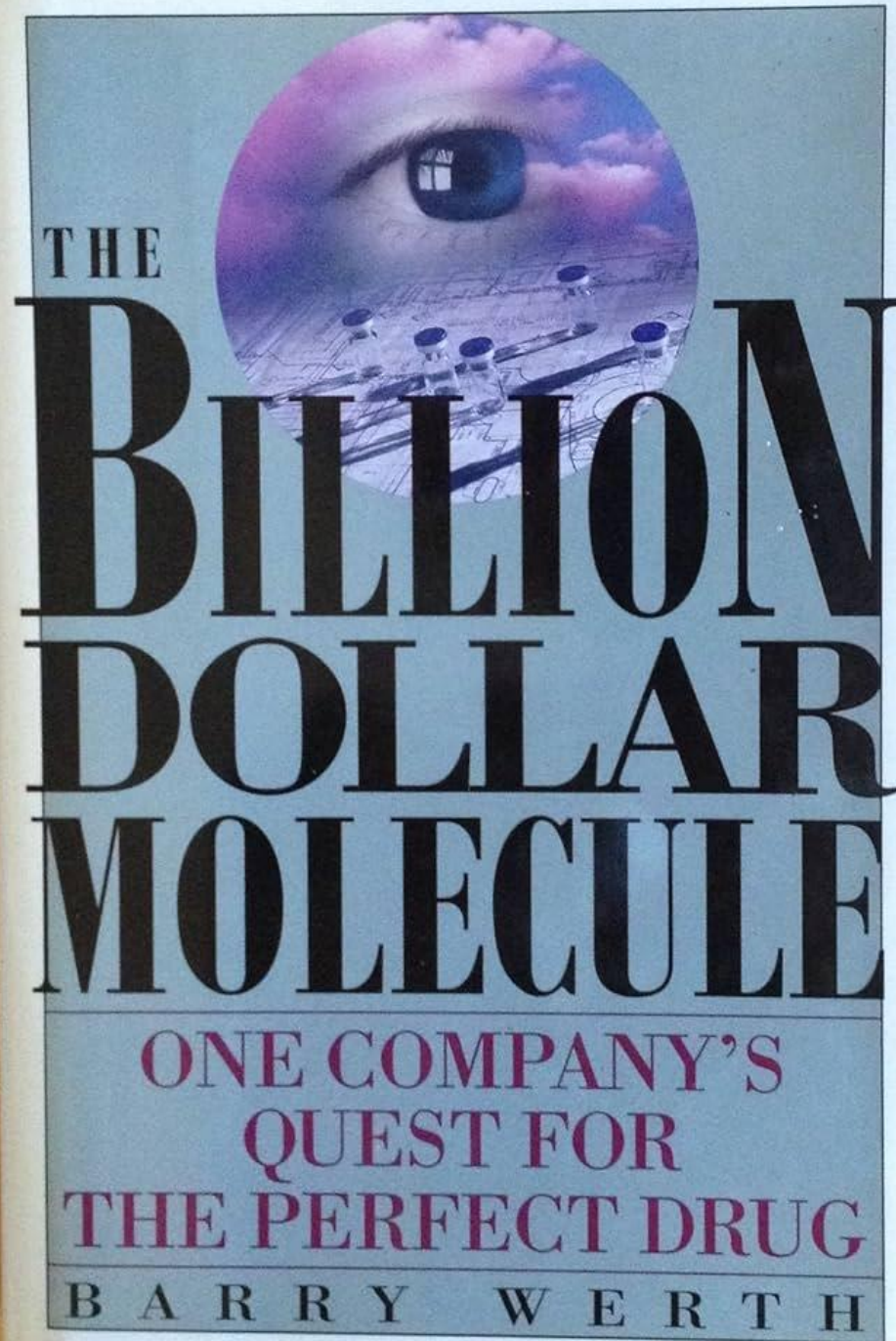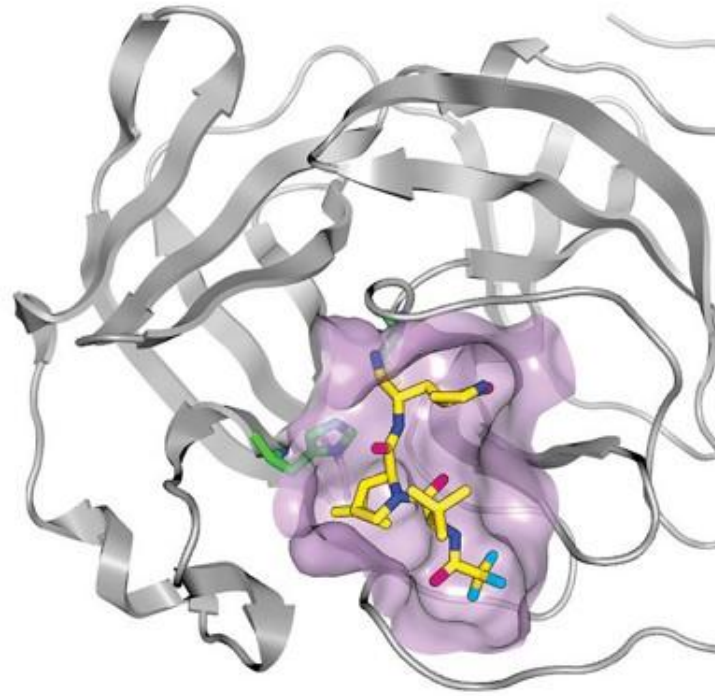
Enoch Huang

15 November 2023

Pfizer

Breakthroughs that change patients' lives
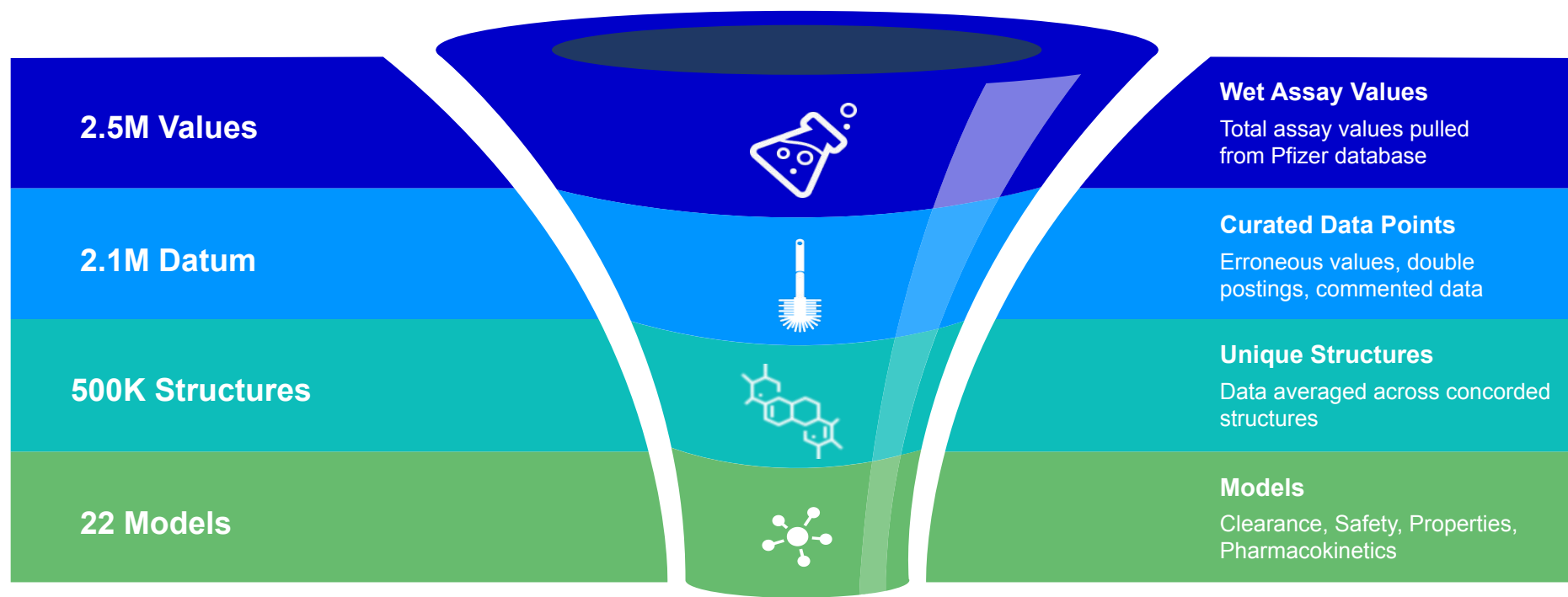
**•Elegance**

**•Reality**

# Machine Learning Tools to Expedite Small Molecule Drug Design

**2.5M Values**

**2.1M Datum**

**500K Structures**

**22 Models**

**Wet Assay Values**
Total assay values pulled from Pfizer database

**Curated Data Points**
Erroneous values, double postings, commented data

**Unique Structures**
Data averaged across concorded structures

**Models**
Clearance, Safety, Properties, Pharmacokinetics

**Models are rebuilt to include new lab data every 2 weeks 🠚 ~80 % of predictions within 2-fold**

The Use of Random Forests for Modeling a Variety of in vitro ADMET Endpoints

A framework for molecular property/activity prediction consisting of a Random Forest model coupled with a custom set of descriptors has been found to be very effective across a variety of endpoints, including kinetic solubility, membrane permeability, metabolic stability, and dofetilide binding.  Random Forests[1] are bagged decision tree ensembles that are trained and applied normally but for one exception: only a small, randomly selected subset of descriptors are considered when selecting the best split at each node during tree construction.  The descriptors used here are all simple molecular substructure or feature counts encoded as Daylight SMARTS queries.  Some mathematical properties of these RF-based models have been explored, including the impact of descriptor and training set selection schemes, nearest neighbor effects, etc.  Additionally, examples will be given to demonstrate that the effectiveness of this modeling paradigm compares favorably to a selection of alternatives.

[1] Breiman, Leo. http://oz.berkeley.edu/users/breiman/.

4

# Our Global ADME/T Machine Learning Models are used ~ 6M / day

| Model | Row Count | % Within 2-Fold |
|---|---|---|
| HHEP Clearance | 92,944 | 74 |
| HLM Clearance | 393,826 | 76 |
| RLM Clearance | 118,201 | 75 |
| RRCK (Pass. Perm.) | 265,074 | 77 |
| NIH MDR (Pgp) ER | 32,598 | 78 |
| BCRP ER | 28,992 | 80 |
| Fu, microsomes | 7,845 | 87 |
| Human Fu, plasma | 10,215 | 69 |
| Rat Fu, plasma | 8,030 | 69 |
| Mouse Fu, plasma | 4,013 | 67 |
| Brain Fu, tissue | 3,216 | 70 |
| Human Blood/Plasma | 2,948 | $R^2$=0.45 |
| Rat Blood/Plasma | 1,742 | $R^2$=0.68 |
| Human Vdss | 1,271 | 62 |
| Rat Vdss | 2,341 | 61 |
| SFLogD | 212,234 | $R^2$=0.78 |
| ELogD | 83,277 | $R^2$=0.86 |
| Kinetic Solubility | 82,996 | 64 |
| Dofetilide Ki | 224,486 | 66 |
| Herg IC50 | 12,963 | 60 |
| THLE IC50 | 101,201 | 77 |
| OATP1B1 Inh | 11,450 | $R^2$=0.67 |

- Design idea prioritization
- Monomer selection in Parallel Medicinal Chemistry (PMC)
- Calculation of PK and Dose

# Most of the effort for a new model is the curation of input data
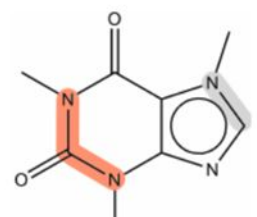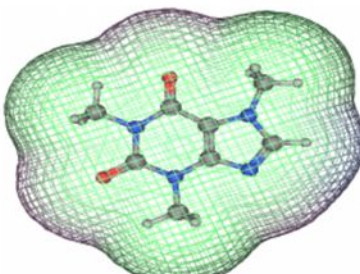
- What assay data is available in the database?
- Are the data suitable?
  - Replicate variability
  - Comment Fields
  - Posting errors
  - Unit errors
- If there are different assays for the same endpoint, can they be combined?
  - Normalization of units of measurement
  - Overlap
  - Correlation
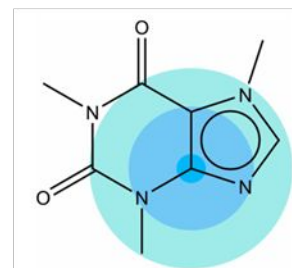- Is the assay updated regularly with new data?

**Pfizer**

**Worldwide Research, Development and Medical**

# 'Classical' machine learning methods: tree-based with descriptors



- XGBoost
- Cubist



| HHEP Metrics | Cubist | DNN | XGBoost |
|---|---|---|---|
| %within 2-Fold | 64 | 60 | 63 |
| Mean Fold Error | 2.2 | 2.3 | 2.1 |
| Pearson's R | 0.56 | 0.50 | 0.59 |
| Spearman's R | 0.66 | 0.61 | 0.70 |

**Worldwide Research, Development and Medical**

# Confidence metrics significantly increased adoption of *in silico* models

- We generate an interpretable probability-based confidence metric

- The score is calibrated via cross-validation to a confidence metric that represents an expected error probability

- The confidence metric captures how close the test compound is to its nearest neighbors in both descriptor space and activity space



Prospective Confidence Metric Performance

$$wRMSD = \sqrt{\frac{\sum_{i=1}^{N} w_i^2 (\hat{y} - y_i)^2}{\sum_{i=1}^{N} w_i^2}}$$

$\hat{y} = $ *predicted value of test compound*

$y_i = $ *actual value of ith neighbor in training set*

$w_i = \dfrac{1}{D + 0.5}$

$D = $ Manhattan Distance between test compound and ith neighbor

Keefer et al (2013) dx.doi.org/10.1021/ci300554t

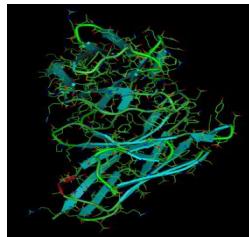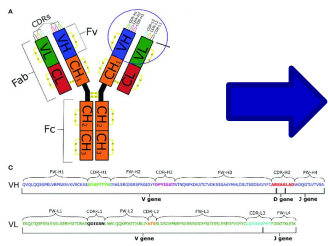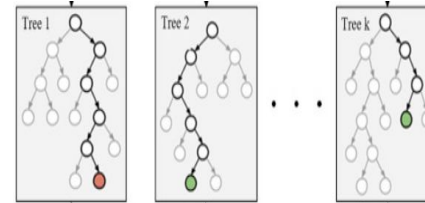# ML models, a computational ecosystem, and culture



## Keys to Success

- Talent, expertise, and remit
- Global, authoritative, standardized data repository
- Infrastructure for publishing, executing, and deploying models
- Confidence scores
- Sophisticated design culture

Pfizer

**Worldwide Research, Development and Medical**

# Attempting the same strategy for large molecules (mAbs) was unsuccessful

Primary sequence

Leave-group-out validation



**Generate features**
Pre-specified set of physchem properties

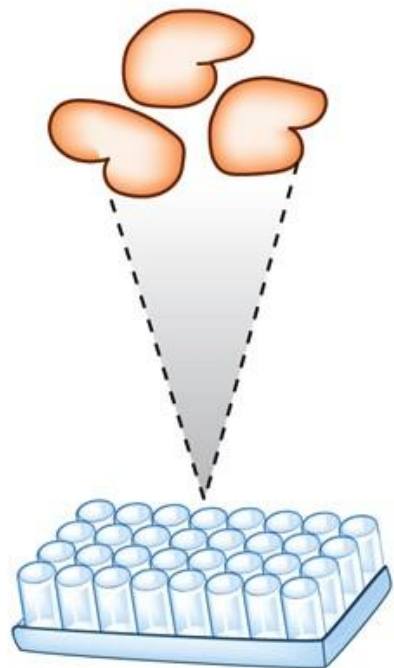**Traditional ML** techniques *failed* to produce generalizable models

Input Layer

6 neurons

100 neurons

500 neurons

200 neurons

50 neurons

Output Layer

Hidden Layers

x1
x2
x3
x4
x5
x6

y1
y2





IMAGENET

# We convert images into lower dimensional vector embeddings



Extract intermediate hidden layer as feature embedding

**Worldwide Research, Development and Medical**
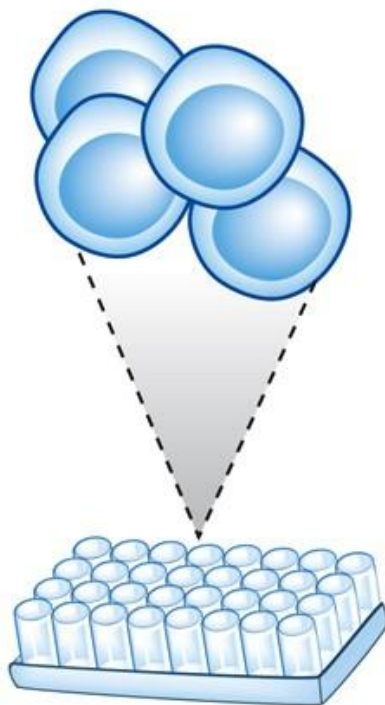
Target-based assays

Cell-based assays
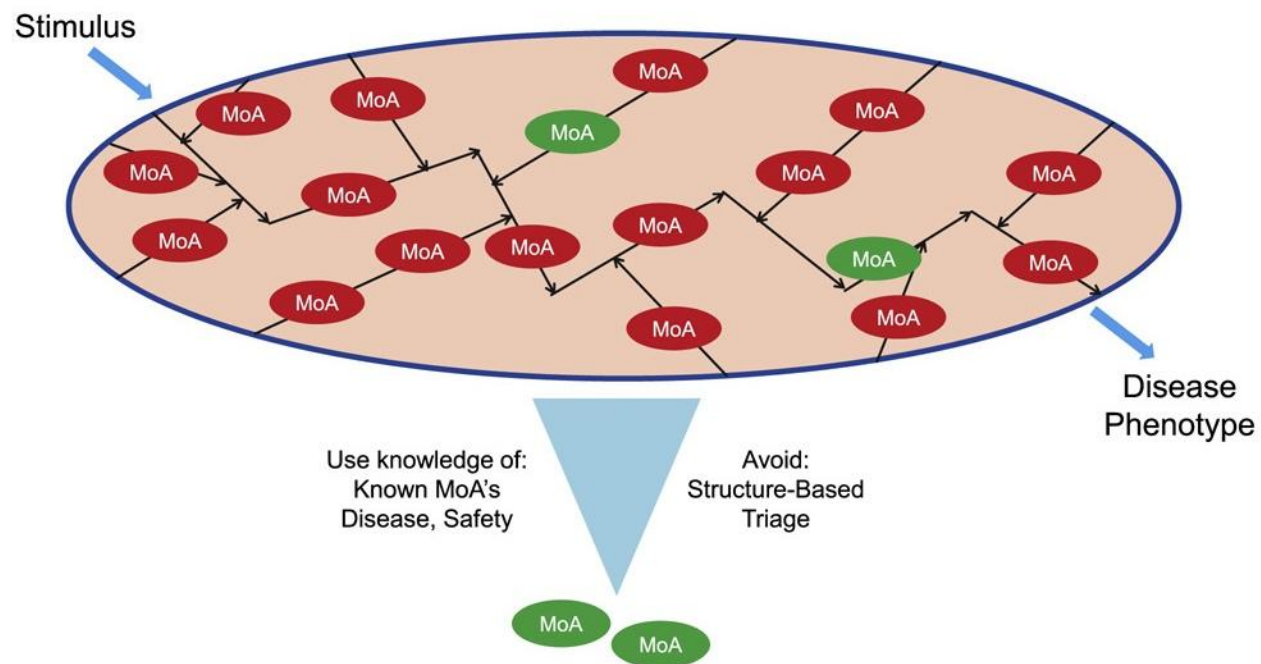
Target-centric
Reductionist view
Validation in cell-based
assays needed

Target-agnostic
Holistic view
More physiologically relevant
Target identification and
validation needed

Stimulus

Disease Phenotype

Use knowledge of:
Known MoA's
Disease, Safety

Avoid:
Structure-Based
Triage

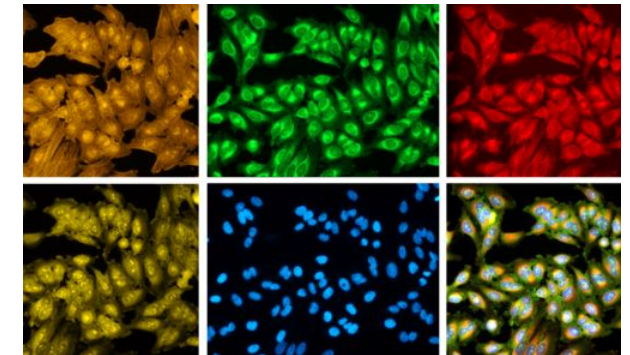Pfizer   **Worldwide Research, Development and Medical**

# Can phenotypic changes of cellular components from compound treatment be learned and associated with specific mechanisms of action (MOAs) via deep learning?
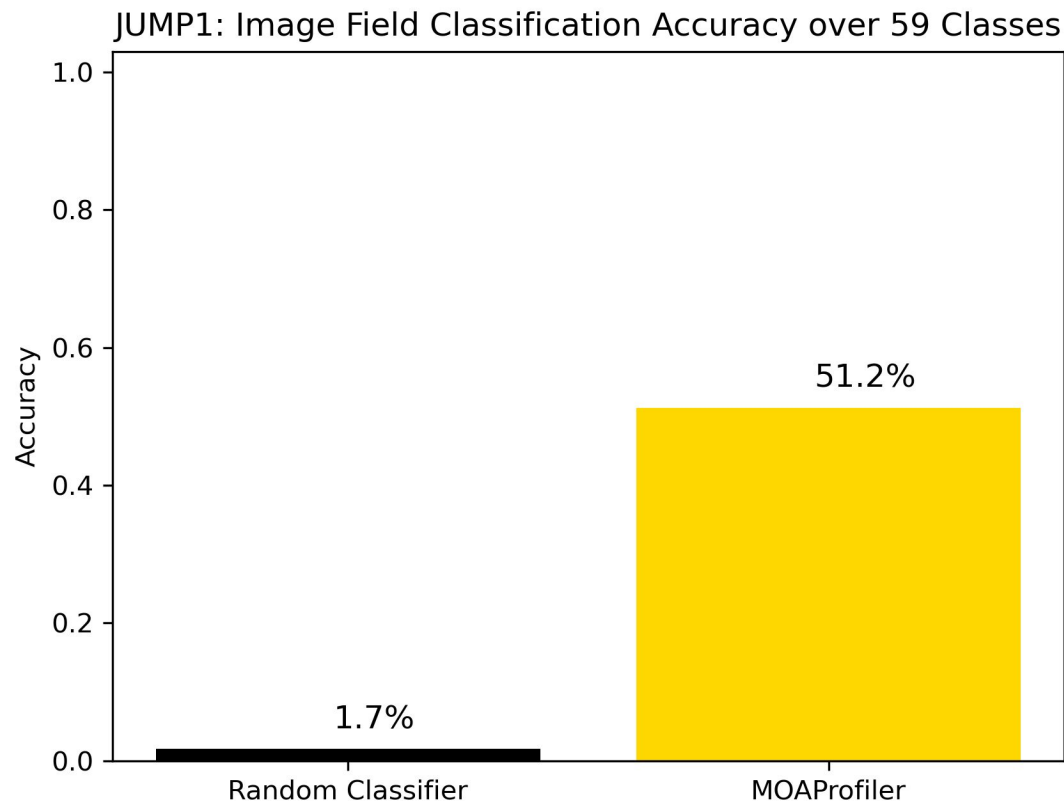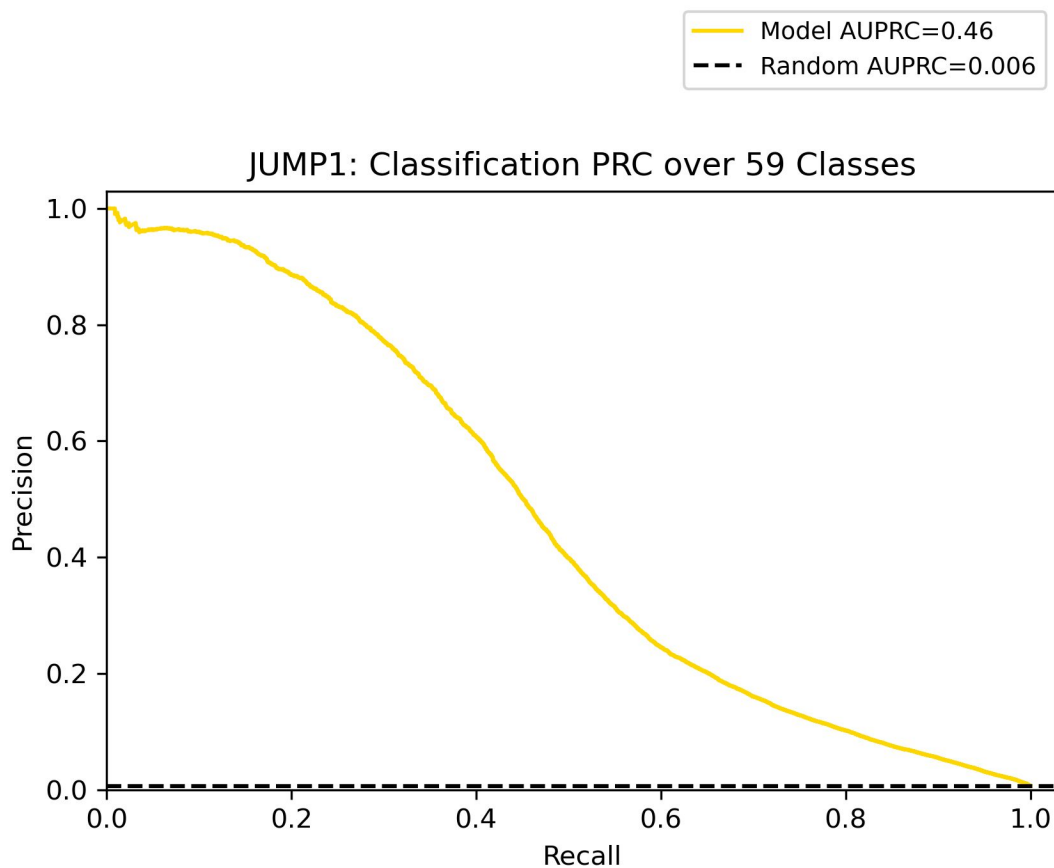
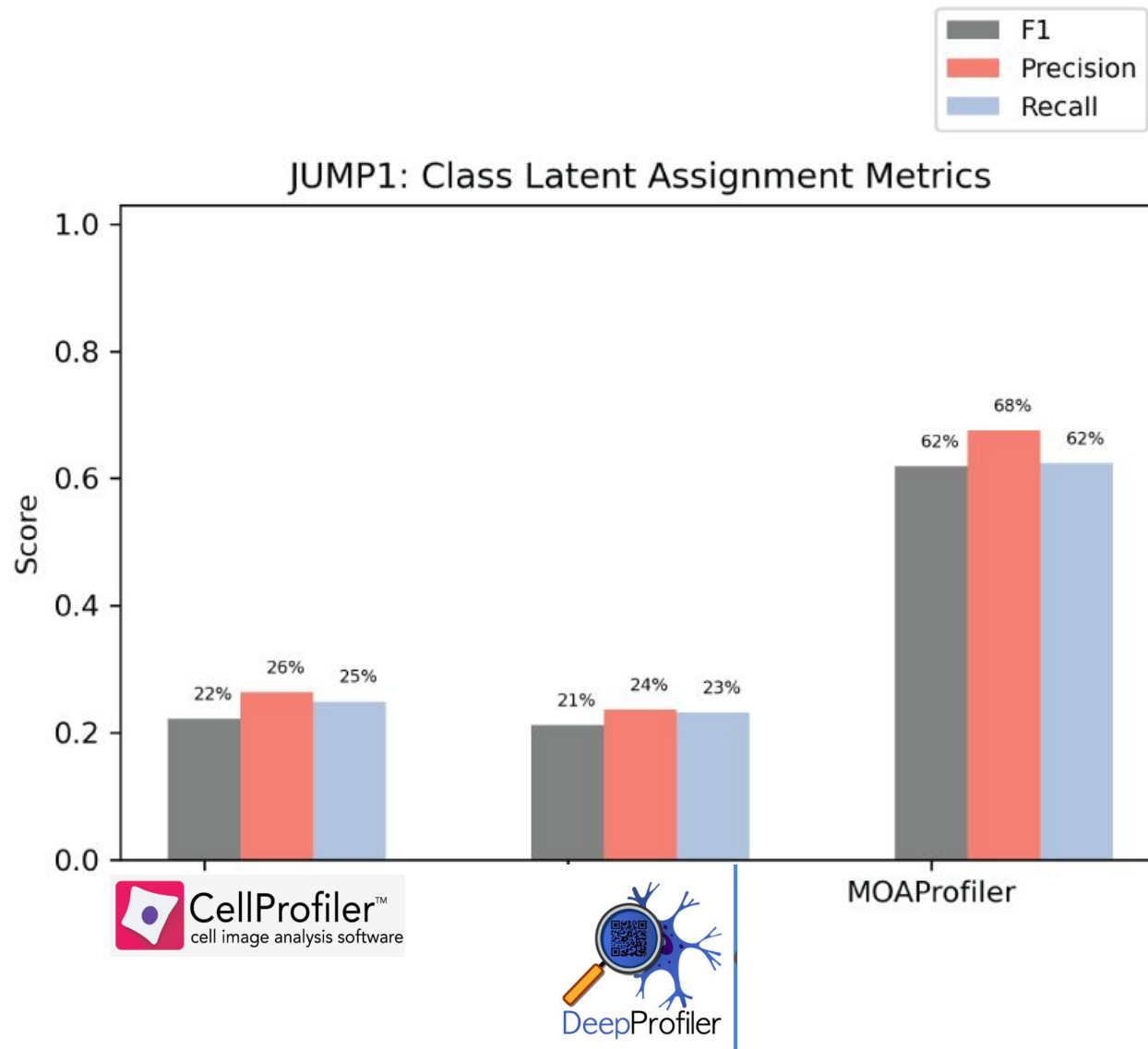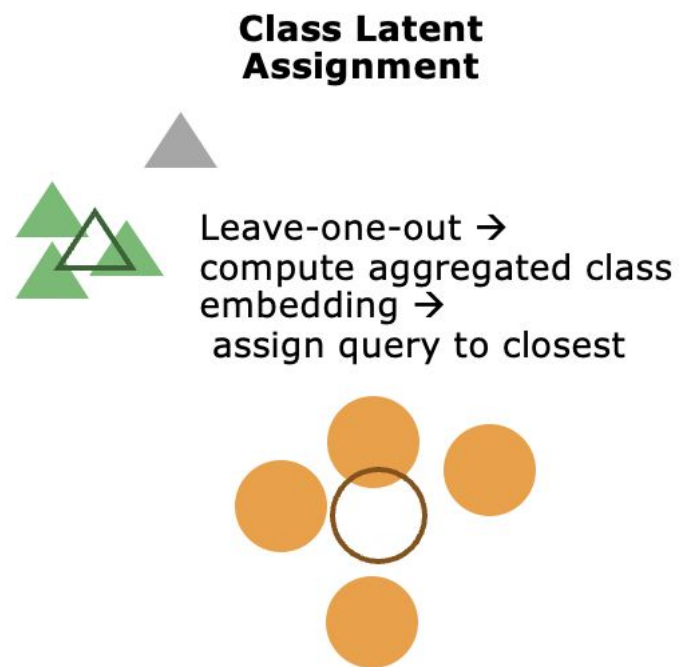- Screen compounds, use computer vision to determine targets/MOAs



Genetic or chemical perturbations → Experiments in multi-well plates → Microscopy imaging → Image analysis → Morphological profiles

Measurements (Features)

Cells

x 384 wells x N plates



- Cell Painting assay (Bray *et al*, 2016)
  - Reveals 8 broadly relevant cellular components or organelles using 6 fluorescent dyes

**Pfizer**

**Worldwide Research, Development and Medical**

# Deep Learning can accurately classify 59 different multi-compound MOAs

# Our embeddings outperform existing methods for MOA class assignment



**Class Latent Assignment**

Leave-one-out → compute aggregated class embedding → assign query to closest

JUMP1: Class Latent Assignment Metrics

Legend:
- F1
- Precision
- Recall

CellProfiler™ — cell image analysis software: F1 22%, Precision 26%, Recall 25%

DeepProfiler: F1 21%, Precision 24%, Recall 23%

MOAProfiler: F1 62%, Precision 68%, Recall 62%

# Case Study: Optimize a Domain of Trispecific Antibody

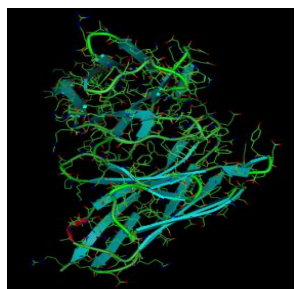*Internally developed AI tool delivers key physical property with speed*
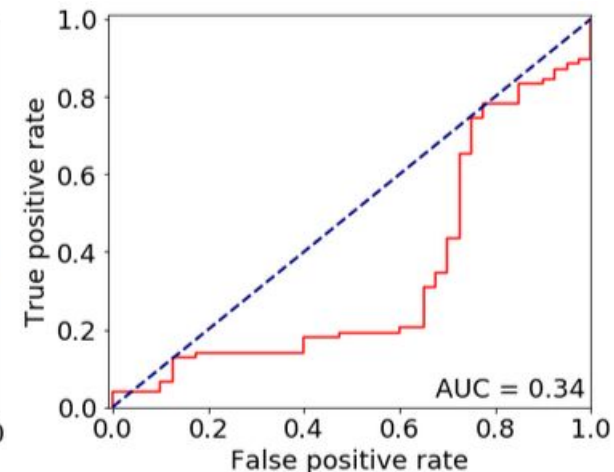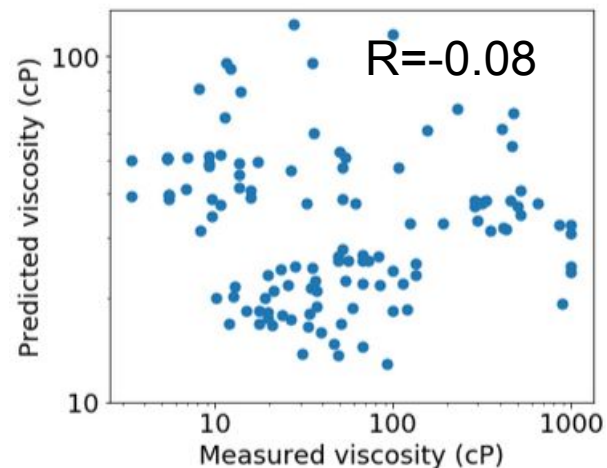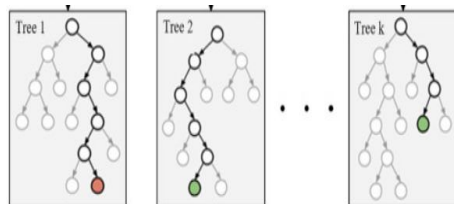


## Situation & Challenges Targeted

- Low antibody viscosity is critical for high dose, low volume subcutaneous delivery and is easier to manufacture

- Traditional viscosity optimization typically requires multiple production / screening cycles

- Scarcity of training data prevented prior AI methods from making accurate antibody viscosity predictions
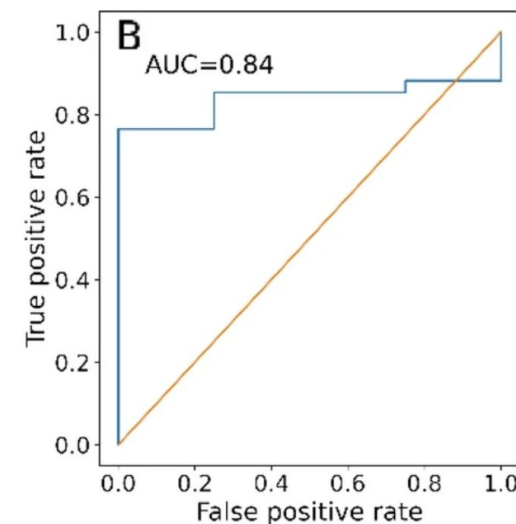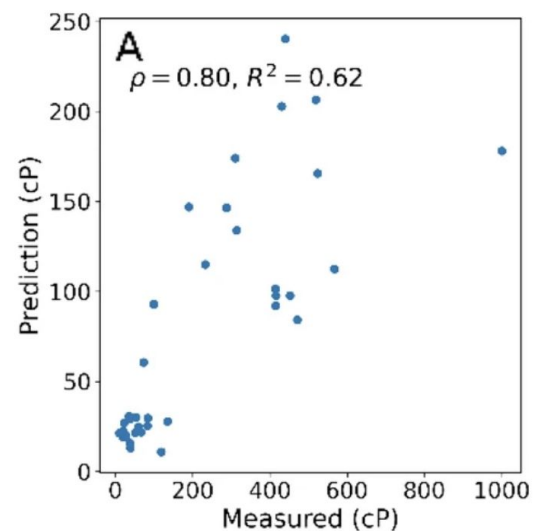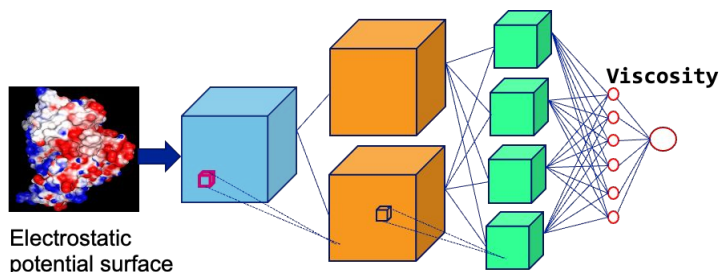
# Using electrostatic potential surface map as the only input to the 3D-CNN prevents overfitting and enables these models to generalize

**Generate features**
Pre-specified set of physchem properties

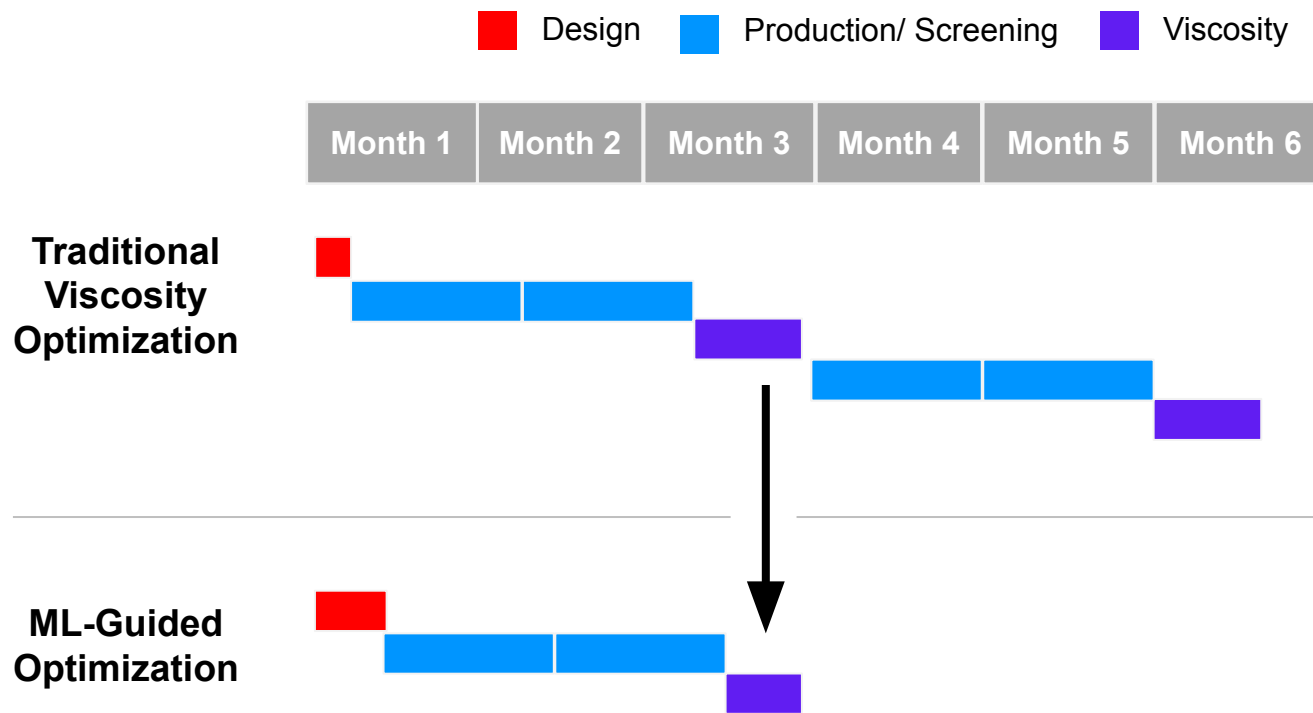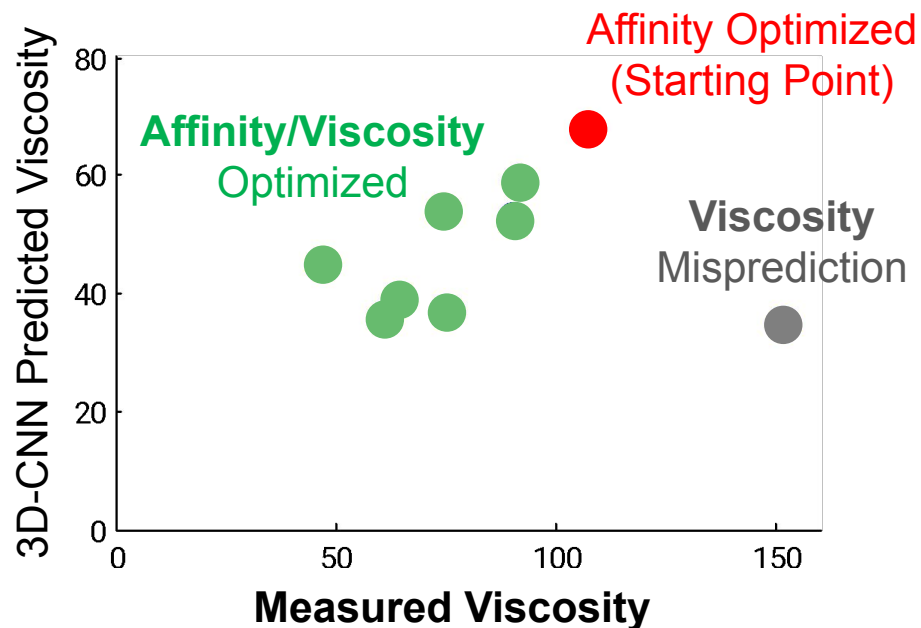**Traditional ML** model failed to generalize



Homology model / crystal structure

Electrostatic potential surface

Viscosity

# ML-Guided Antibody Viscosity Optimization was at Least 50% Faster*

**AI prediction correlated strongly with measured viscosity of optimized mutants**



Affinity Optimized (Starting Point)

Affinity/Viscosity Optimized

Viscosity Misprediction

3D-CNN Predicted Viscosity

Measured Viscosity

Design | Production/ Screening | Viscosity

| | Month 1 | Month 2 | Month 3 | Month 4 | Month 5 | Month 6 |
|---|---|---|---|---|---|---|

Traditional Viscosity Optimization

ML-Guided Optimization

**Prioritization of Antibody Mutants for Testing Eliminated Need for Multiple Production / Screening Cycles**

*Compared to traditional approach

# We are Leveraging Recent Advances in Language Modeling Techniques to Support Large Molecule Discovery Efforts

- Recent advances in AI can be attributed to one methodological breakthrough in deep learning: *Transformers*

## Attention Is All You Need

Google Research          *Published: 2017*

**Ashish Vaswani\***       **Noam Shazeer\***       **Niki Parmar\***       **Jakob Uszkoreit\***
Google Brain             Google Brain             Google Research         Google Research

## Prominent Transformer-based AI models

- Generative Pre-trained Transformer (GPT)   OpenAI ChatGPT

- AlphaFold2

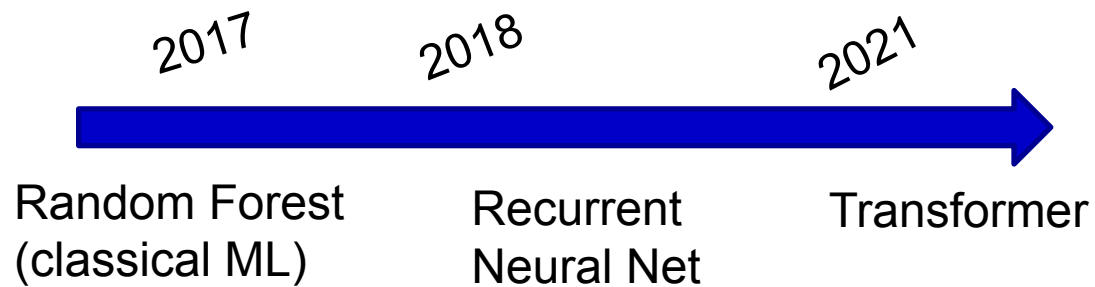We use Transformers in predictive modeling efforts for three tasks:

1. Antibody clearance

2. mRNA design

3. mAb immunogenicity risk assessment

Pfizer

# Non-specificity Predictions from *minGPT*[*] based Models are being used to Reduce mAb PK Risk in the Early Discovery Stage
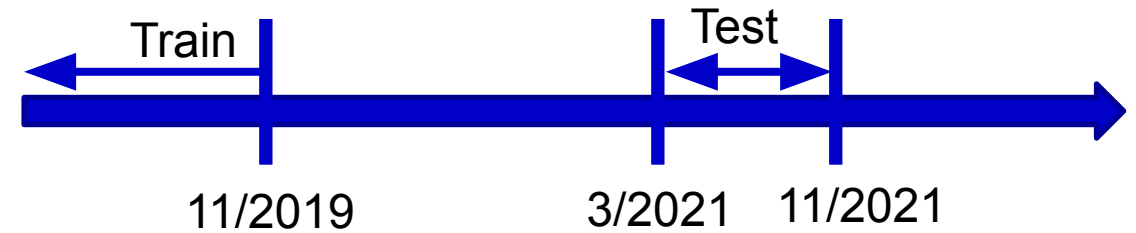
- *In vitro* non-specificity endpoints correlate well with *in vivo* clearance (Avery et al., *mAbs* 2018)
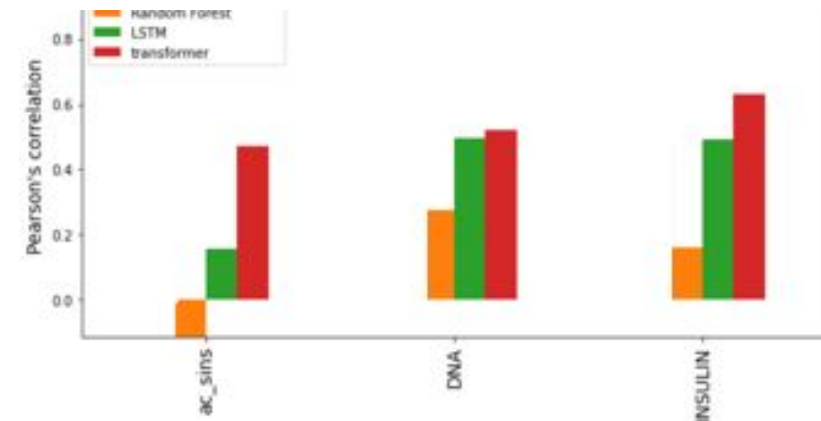


- Our choice of ML techniques over time



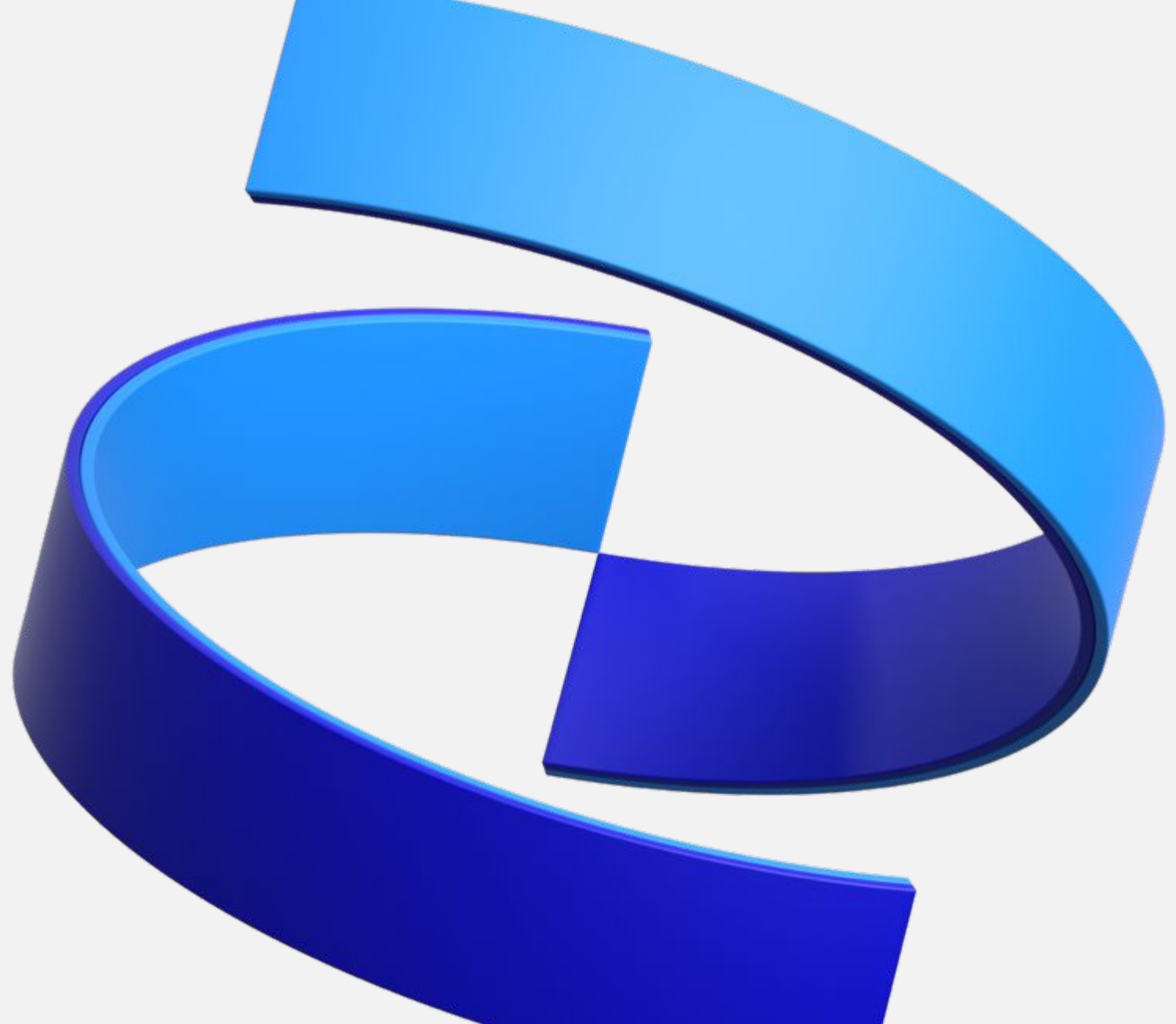- Adoption of advanced ML techniques have led to better prediction performance



Test on antibodies from new portfolio projects



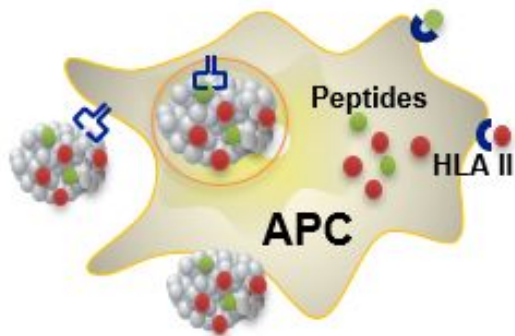- Models have been integrated into the Pfizer developability assessment workflow

# Thanks to my Pfizer colleagues

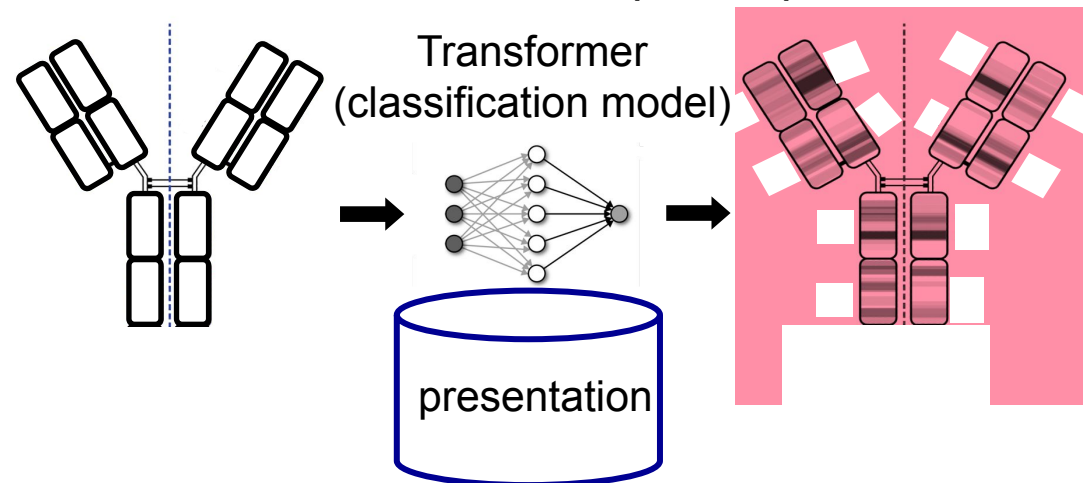Chris Keefer
Daniel Wong
Brajesh Rai

# Transformer Models have been Developed to Identify Potential Epitopes on Antibody Sequences

- Therapeutic antibodies run the risk of being recognized as foreign by a host immune system
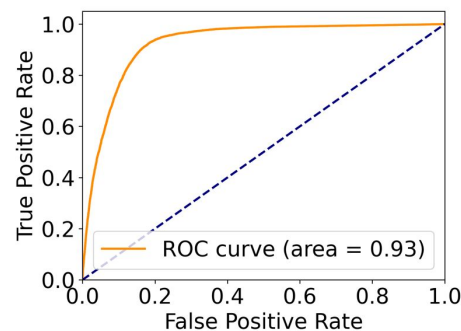


- Current immunogenicity risk assessment relies on peptide-HLA II binding predictions
  - Trained on *in vitro* binding affinities

- Recent publications have shown that peptide-HLA II presentation is a better predictor of immunogenicity
  - Trained on MS immunopeptidomic data
  - Chen, Nat Biotechnol 2019; MARIA, Stanford

Prediction of HLA II allele-specific presentation



Transformer (classification model)

presentation

Retrieval of true HLA2 presented peptides



ROC curve (area = 0.93)

**Epitope prediction accuracy (96 mAbs): 97%**

Binding-based (current): 78%

MIT
Technology
Review

Featured     Topics     Newsletters     Events     Podcasts
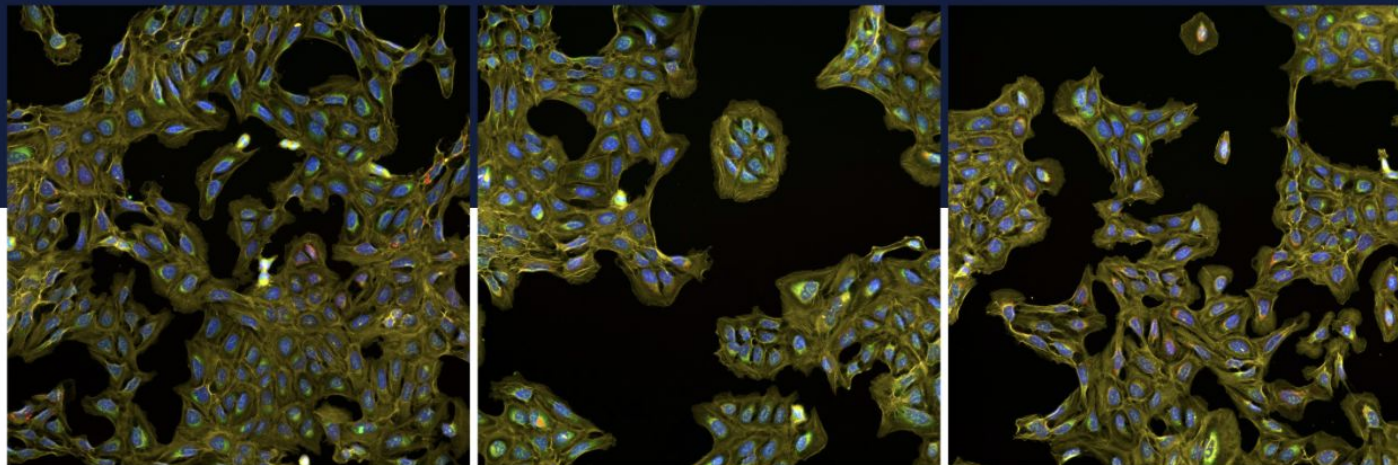
SIGN IN

SUBSCRIBE

BIOTECHNOLOGY AND HEALTH
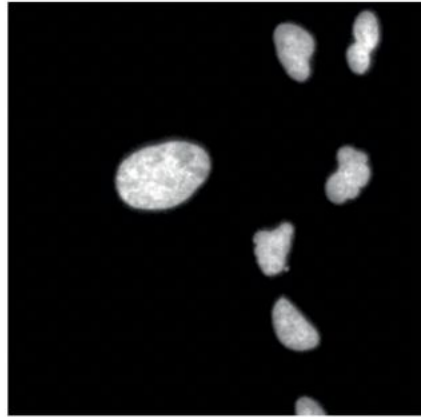
# A technique called Cell Painting could speed drug discovery

A public-private consortium has released a huge collection of image-based cell profiles.
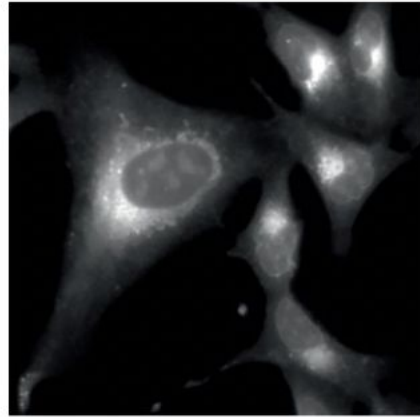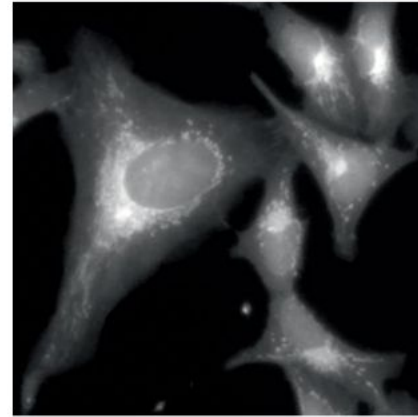
By Esther Landhuis

March 3, 2023



Pfizer

Worldwide Res

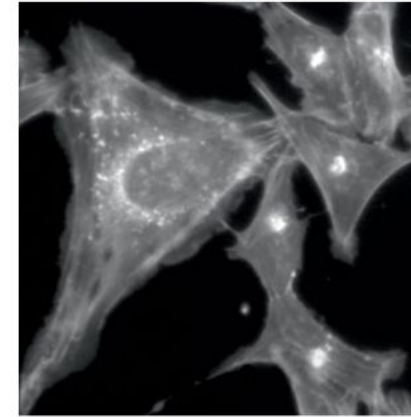# Cell Painting provides phenotypic information about key cellular structures
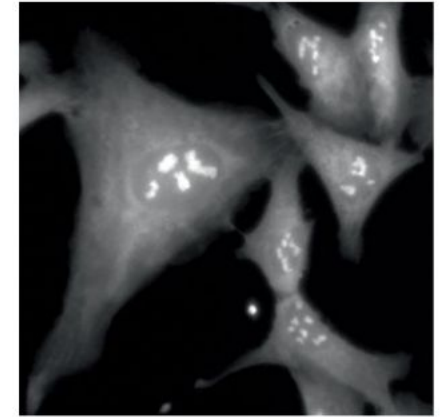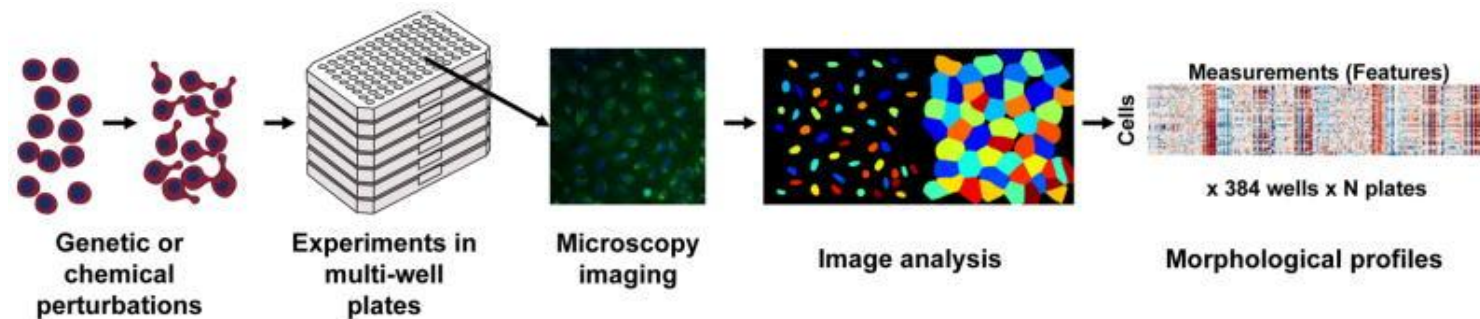


Hoechst 33342
Nucleus

Alexa 488
Endoplasmic Reticulum

Alexa 647
Mitochondria

Alexa 568
Actin, Golgi Apparatus, Plasma Membrane

Alexa long
Cytoplasmic RNA

Genetic or chemical perturbations → Experiments in multi-well plates → Microscopy imaging → Image analysis → Morphological profiles

Measurements (Features)
Cells
x 384 wells x N plates
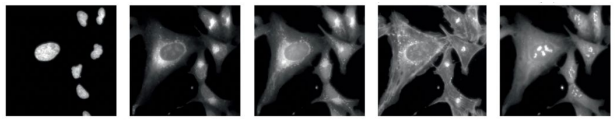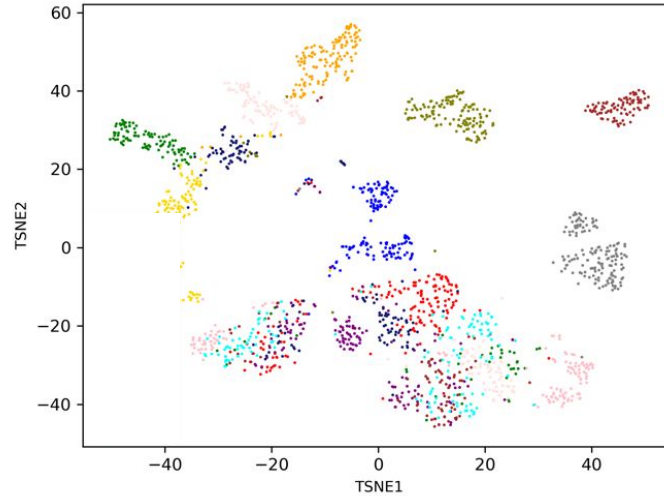
Chandrasekaran et al. Image-based profiling for drug discovery: due for a machine-learning upgrade? Nature Reviews Drug Discovery (2020).
Bray et al. Cell Painting, a high-content image-based assay for morphological profiling using multiplexed fluorescent dyes. *Nat Protocol (*2016).

Pfizer

# We aggregate individual image embeddings into well-level embeddings to compare with CellProfiler and DeepProfiler
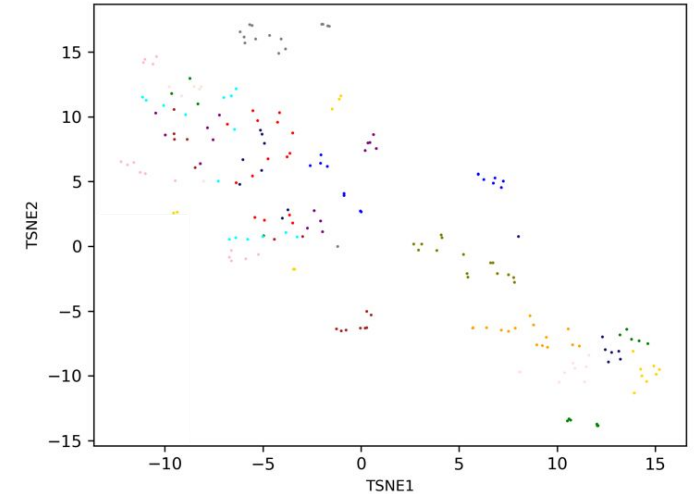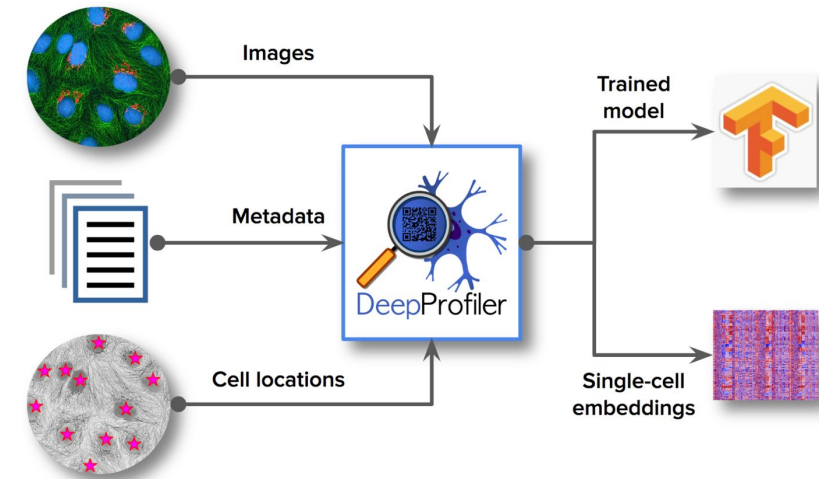


Encode image fields via neural network

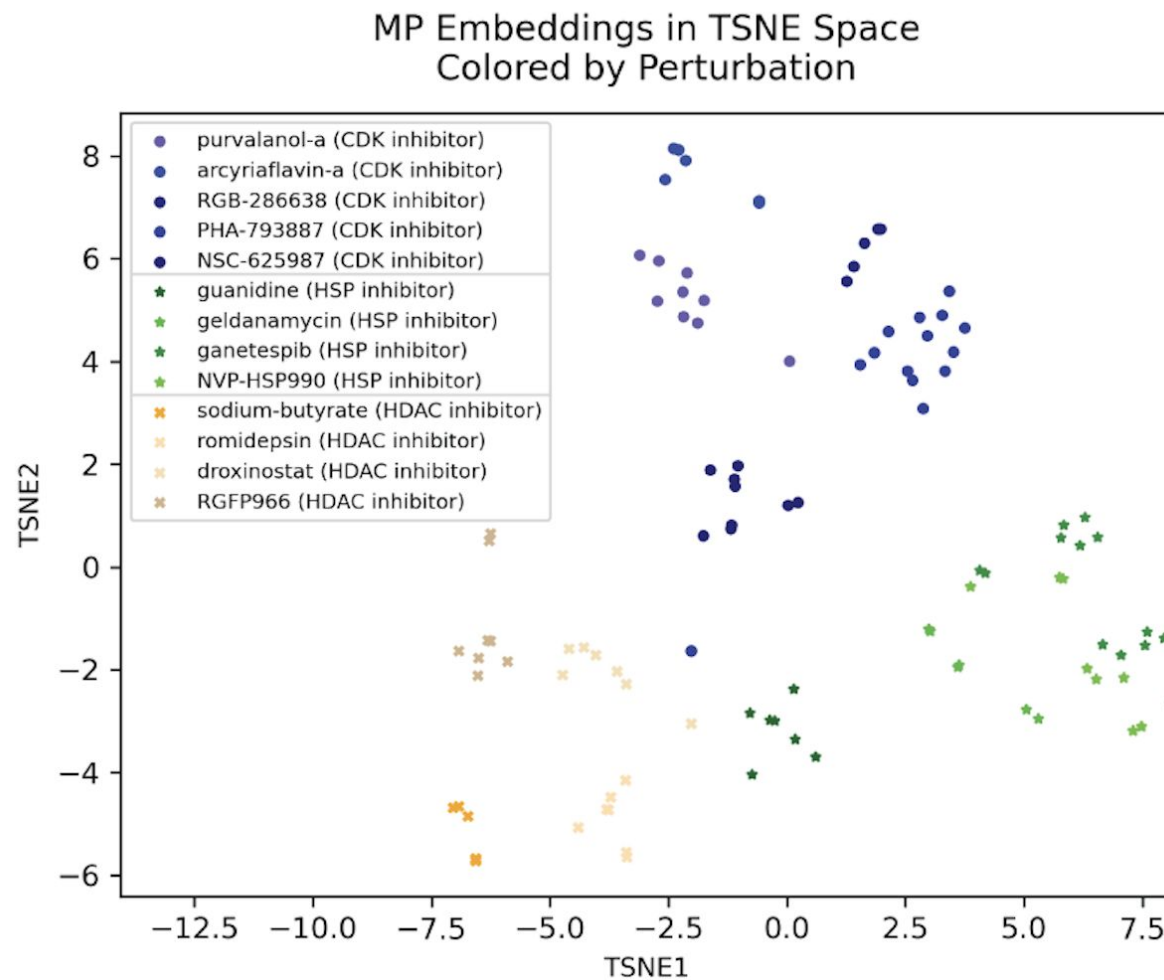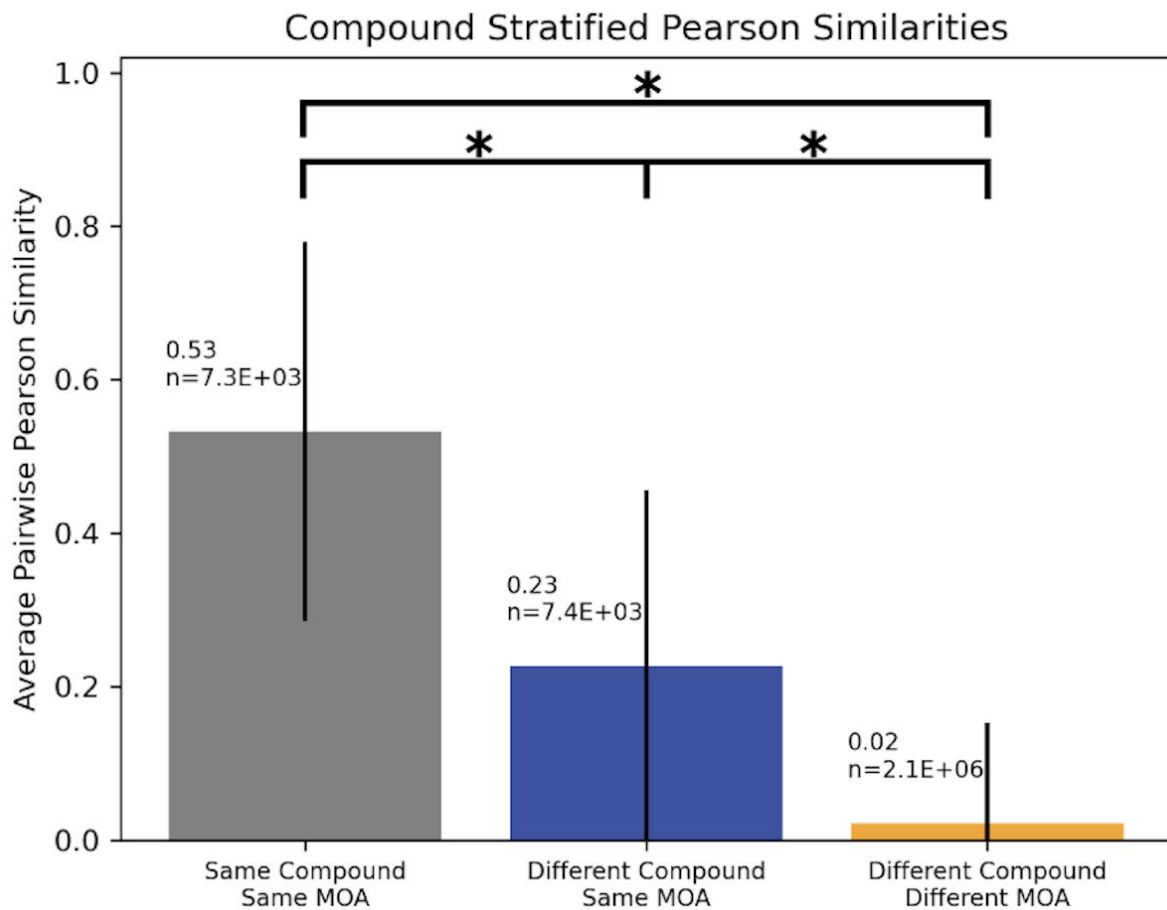Visualization of embeddings at the image level

aggregate into well-level embeddings

Carpenter et al. CellProfiler: image analysis software for identifying and quantifying cell phenotypes. Genome Biology (2006).
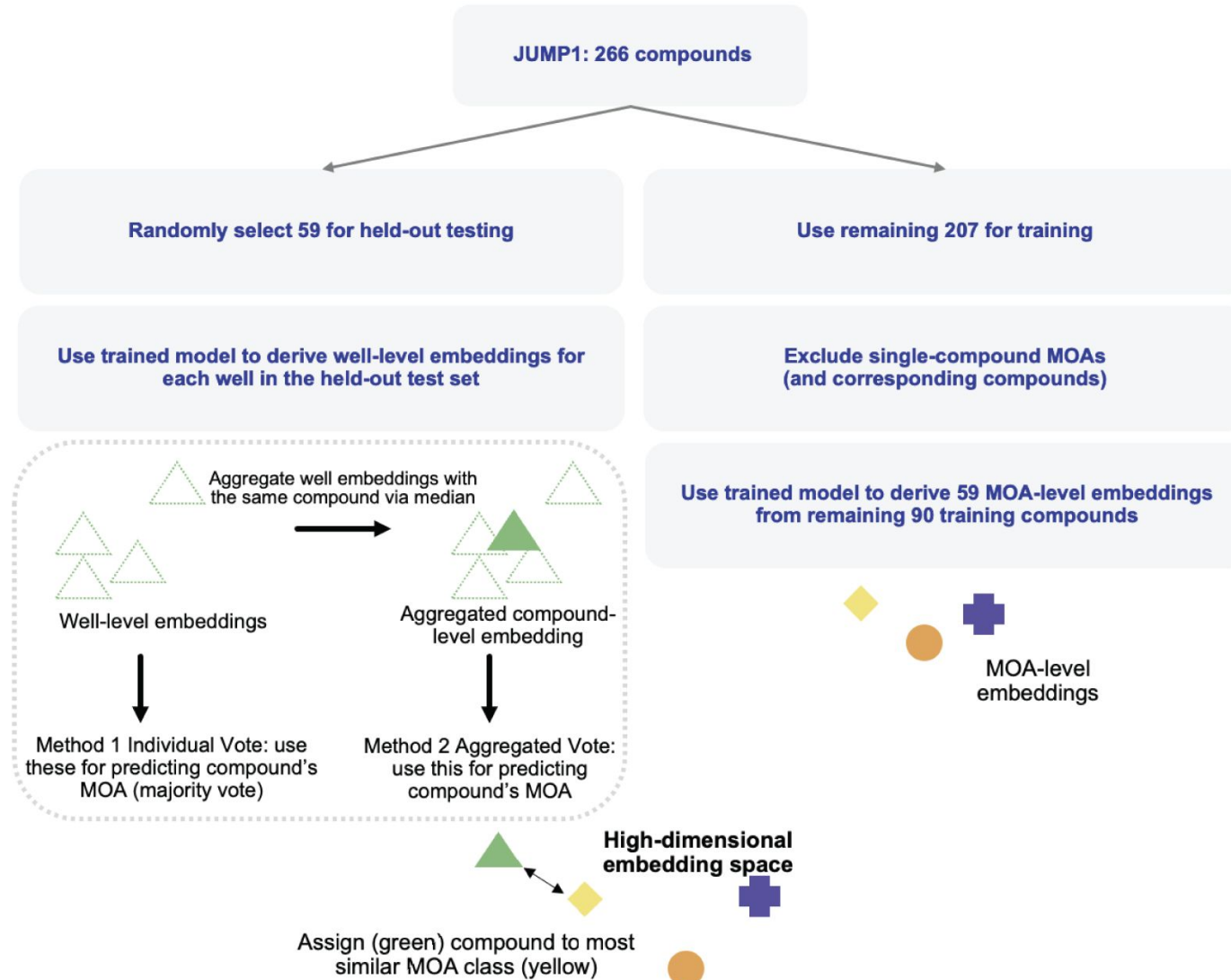https://github.com/cytomining/DeepProfiler

Images

Metadata

Cell locations

DeepProfiler

Trained model

Single-cell embeddings

# Model learns MOA-specific phenotypes amidst a diverse compound space



Compound Stratified Pearson Similarities

MP Embeddings in TSNE Space
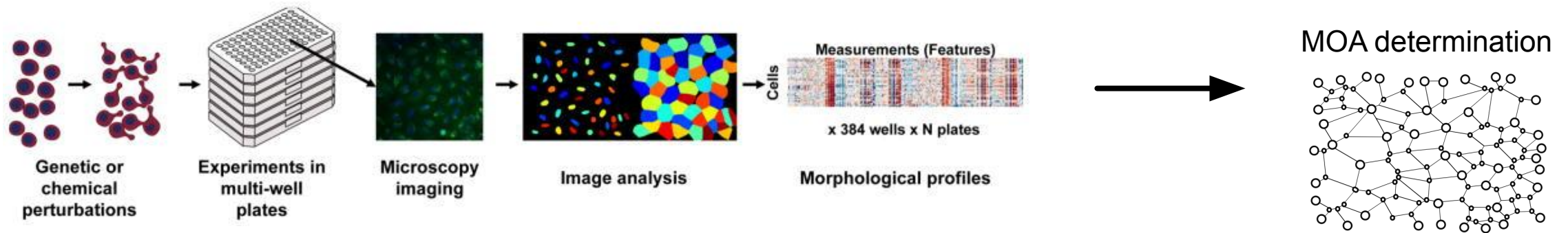Colored by Perturbation

# Can we predict the MOA of held-out compounds never exposed to model training?

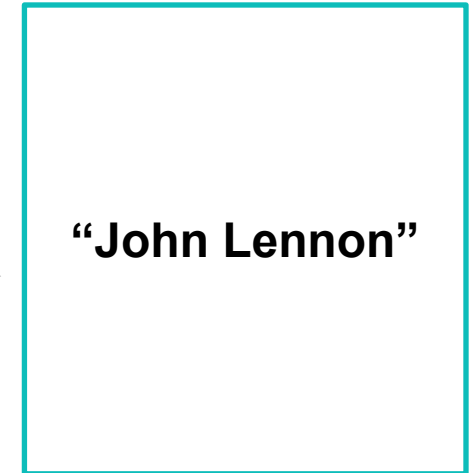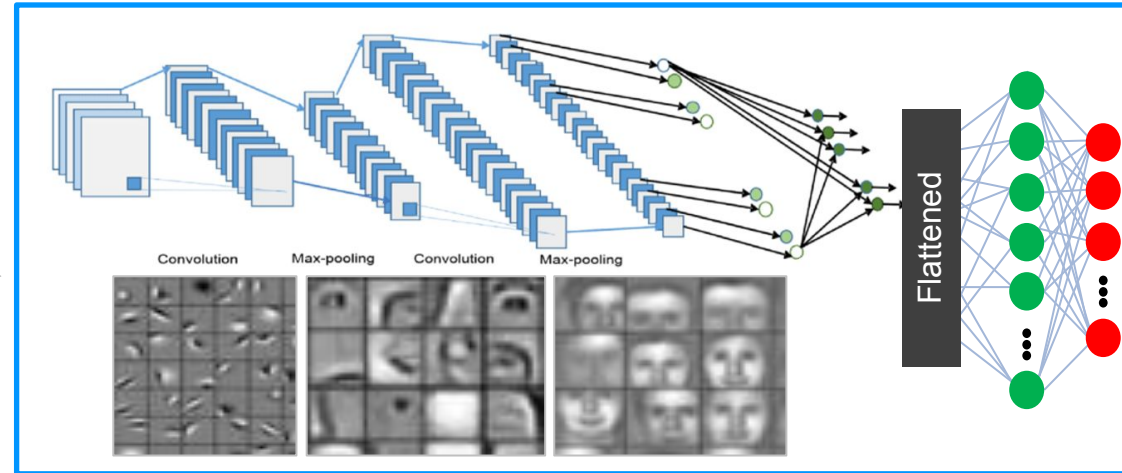Same set up but hold out *compounds* instead of wells!

# Conclusions

- Drug MOA determination via CellPainting phenotypes is possible!

- MOAProfiler is more performant for MOA determination than both the gold-standard CellProfiler and DeepProfiler

- Approach generalizes to two different datasets and predicts the MOA of *held-out* compounds
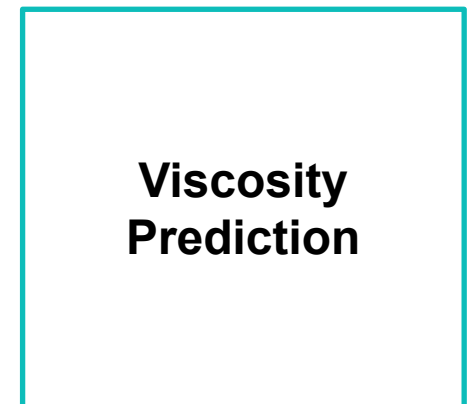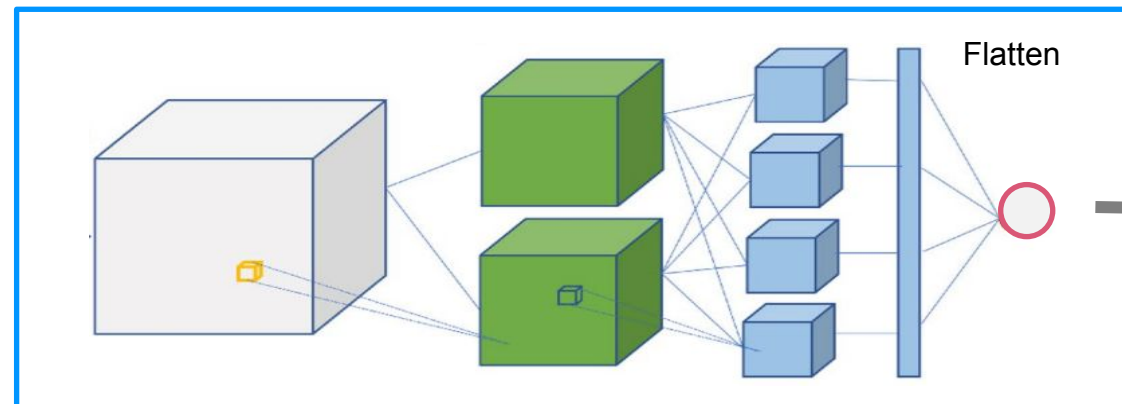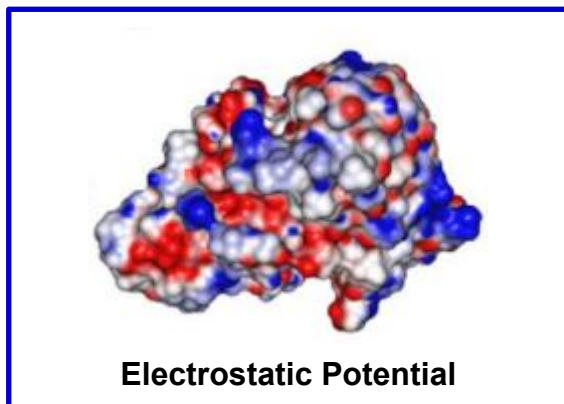
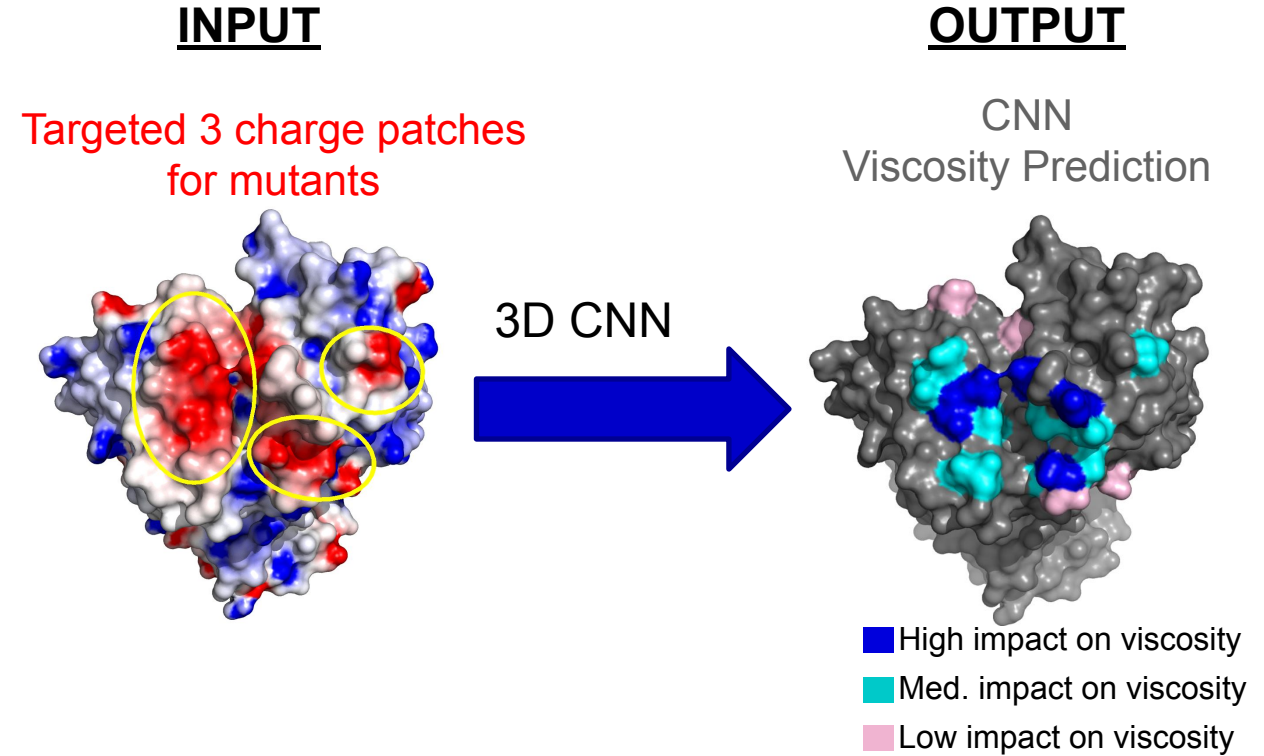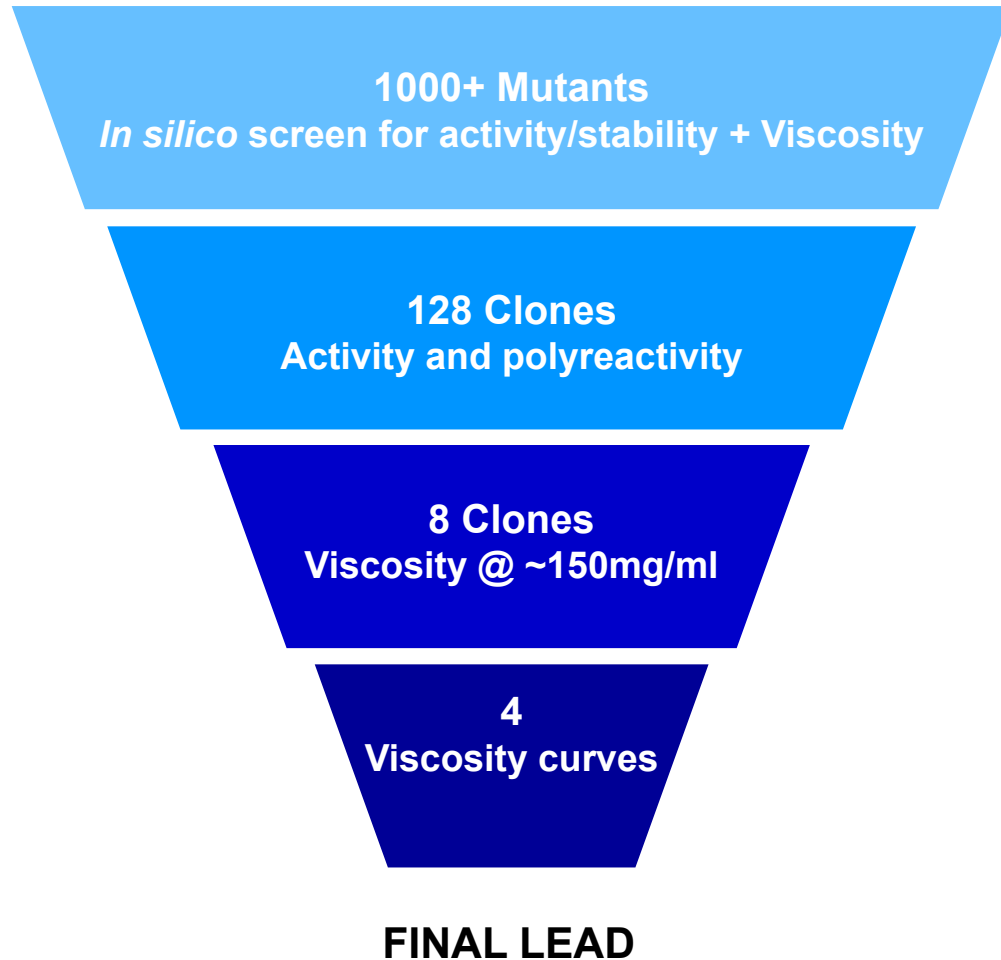# PfAbNet: Pfizer's internally developed 3D CNN that predicts antibody viscosity



**Traditional 2D CNN**

INPUT — FEATURE EXTRACTION — OUTPUT

"John Lennon"

**PfAbNet**

Electrostatic Potential

Flatten

Viscosity Prediction

Rai et al. (2023) https://doi.org/10.1038/s41598-023-28841-4

# AI Guided Optimization Delivering Antibodies in Less than Half the Time



**1000+ Mutants**
*In silico* screen for activity/stability + Viscosity

**128 Clones**
Activity and polyreactivity

**8 Clones**
Viscosity @ ~150mg/ml

**4**
Viscosity curves

**FINAL LEAD**

**INPUT**

Targeted 3 charge patches
for mutants

3D CNN

**OUTPUT**

CNN
Viscosity Prediction

- High impact on viscosity
- Med. impact on viscosity
- Low impact on viscosity

# 3D-CNN model outperforms previous methods[14,15] in viscosity prediction

From: Low-data interpretable deep learning prediction of antibody viscosity using a biophysically meaningful representation

Article | Open Access | Published: 20 February 2023

## Low-data interpretable deep learning prediction of antibody viscosity using a biophysically meaningful representation
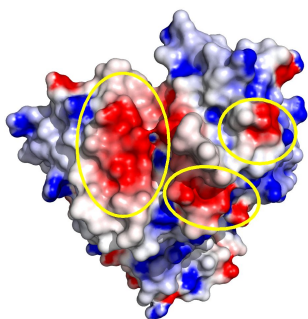
Brajesh K. Rai ✉, James R. Apgar & Eric M. Bennett

14. Agrawal, N. J. *et al.* Computational tool for the early screening of monoclonal antibodies for their viscosities. *MAbs* **8**, 43–48. https://doi.org/10.1080/19420862.2015.1099773 (2016).

15. Sharma, V. K. *et al.* In silico selection of therapeutic antibodies for development: Viscosity, clearance, and chemical stability. *Proc. Natl. Acad. Sci.* **111**, 18601–18606 (2014).

# Application of PfAbNet to a Trispecific Antibody

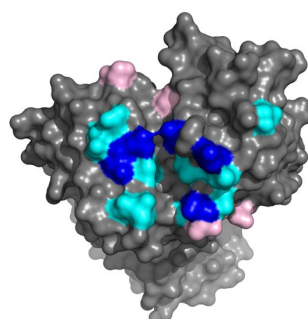*AI prediction correlated strongly with measured viscosity of optimized mutants*

**INPUT**

Targeted 3 charge patches for mutants

**OUTPUT**

PfAbNet Viscosity Prediction

■ High impact on viscosity
■ Med. impact on viscosity
■ Low impact on viscosity

**VALIDATION:**

One round of design delivered improved viscosity

Pfizer Worldwide Research, Development and Medical