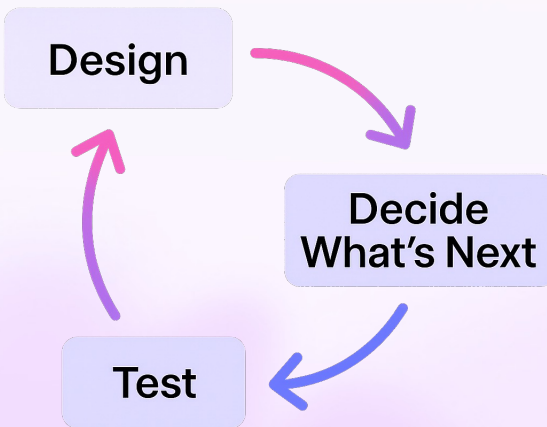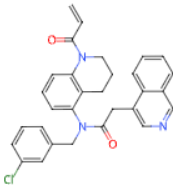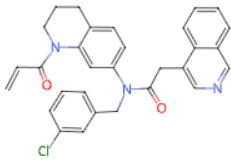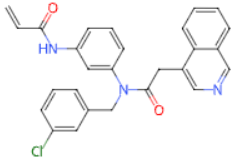# Bio-Assay Derived Fingerprints: Performance Comparison in AutoML Pipelines
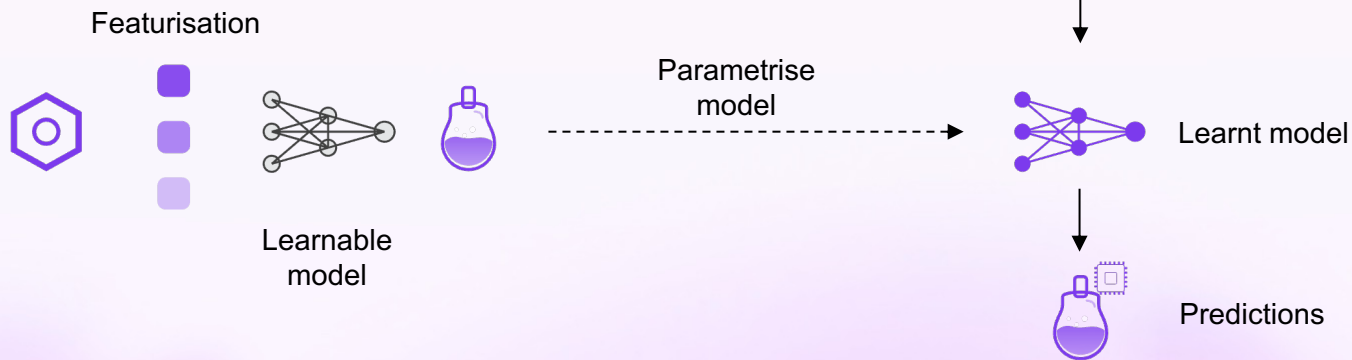
Ryan Greenhalgh | CTO and Co-Founder

# Early Discovery: What to Make Next?

- Early programs <100 datapoints per assay
- Need to **prioritise next compounds**
- Can we build models to help?

Design → Decide What's Next → Test → Design



| | Molecules | pIC50 | LogD | MLM (uL/min/mg) |
|---|---|---|---|---|
| 1 | | 6.32 | 3.3 | 1460 |
| 2 | | None | None | 524 |
| 3 | | 5.57 | 3.8 | None |

deepmirror

Program data = $\left[ \begin{array}{cc} \text{⬡} & \text{🧪} \end{array} \right]$

Featurisation

Learnable model

Parametrise model

New ideas

Learnt model

Predictions

deepmirror

Polaris | ASAP | Open Molecular Software Foundation

# Antiviral Competition

Ligand Poses  Ligand Potency  Ligand ADMET

LABORATORY  CLINIC

MULTIPLE MOLECULES
ABUNDANCE OF DATA

SINGLE DRUG CANDIDATE
SCARCITY OF DATA

# MolE: a foundation model for molecular graphs using disentangled attention

Oscar Méndez-Lucio [1], Christos A. Nicolaou [1,2] & Berton Earnshaw [1]

# How to learn on fingerprints

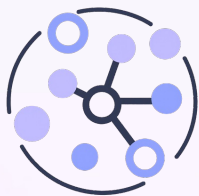| Model | Pros | Cons | <100 datapoints | 100–10k | >10k |
|---|---|---|---|---|---|
| Gaussian Process | Uncertainty-aware, good w/ low data | Slow, tricky to scale | +++ | +++ | ! |
| Trees / SVM / kNN | Fast, robust | No uncertainty, less expressive | ++ | +++ | ! |
| Deep Learning | Best w/ big data | Overfits on small datasets | xx | ! | +++ |

deepmirror

C=CC(=O)N1CCCc2...

Cheminformatics

Mordred

RDKit

Morgan

MACCS

fast, proven, interpretable

Machine Learning

CheMeleon

ChemBERTa

Mol2Vec

# RDKit / Mordred / MACCS

Physicochem    Topology    Electronic    Fragments

Descriptor value

low                                                    high

Does the molecule have 5M ring?

ON (1)        OFF (0)

# Mol2vec/ChemBERTa

Leveraging a CBOW/transformer trained on large datasets i.e. PubChem dataset (77M compounds).
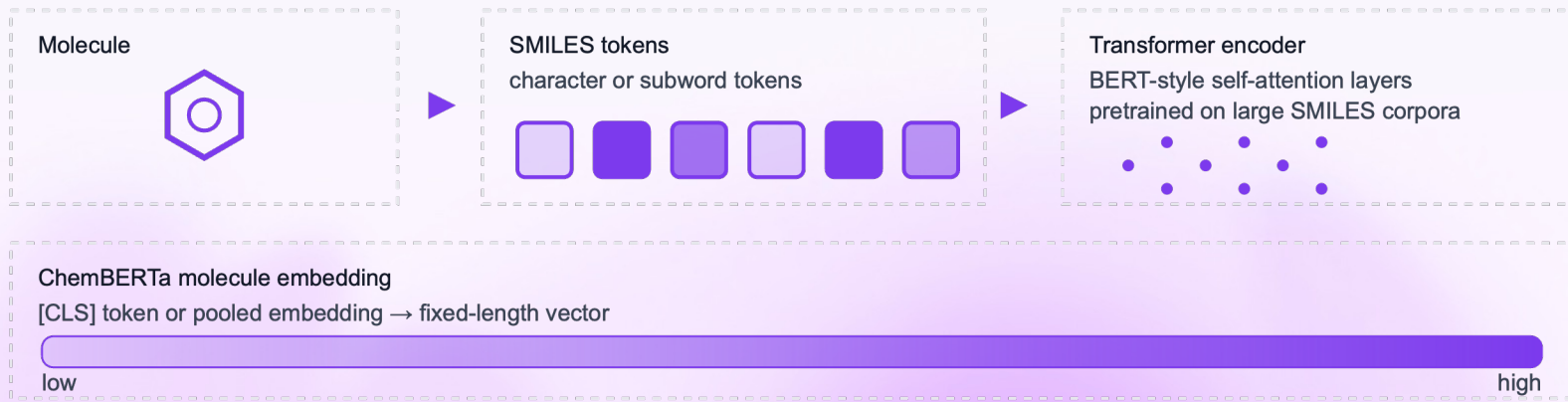
Seyone Chithrananda
University of Toronto
seyone.chithrananda@utoronto.ca

Gabriel Grand
Reverie Labs
gabe@reverielabs.com

Bharath Ramsundar
DeepChem
bharath.ramsundar@gmail.com

**Mol2vec: Unsupervised Machine Learning Approach with Chemical Intuition**

Sabrina Jaeger, Simone Fulle*, and Samo Turk*

Molecule

SMILES tokens
character or subword tokens

Transformer encoder
BERT-style self-attention layers
pretrained on large SMILES corpora

ChemBERTa molecule embedding
[CLS] token or pooled embedding → fixed-length vector
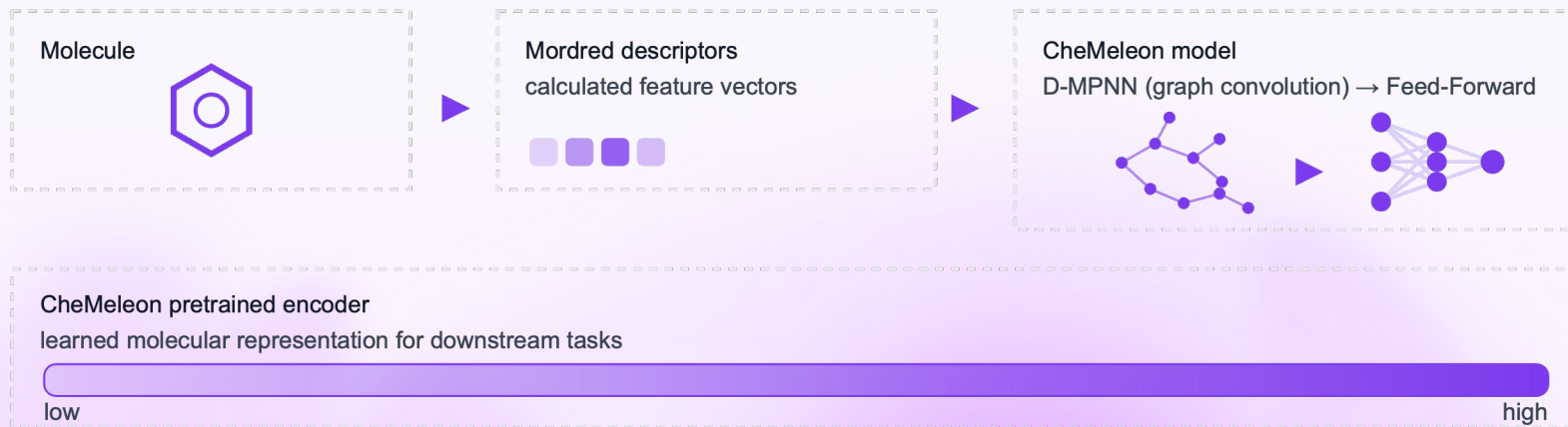
low                                                                                          high

# CheMeleon

Similar to ChemBERTa, but
applied to physical fingerprints
and implemented with GNNs.



**Molecule**

**Mordred descriptors**
calculated feature vectors

**CheMeleon model**
D-MPNN (graph convolution) → Feed-Forward

**CheMeleon pretrained encoder**
learned molecular representation for downstream tasks

low                                                                high
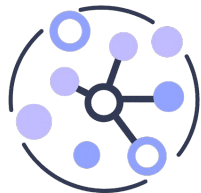
**deepmirror**

CheMeleon

ChemBERTa

Mordred

Mol2Vec

Morgan

RDKit

MACCS

Traditional fingerprints rely on physico-chemical properties, structural chemistry, or information-theoretic descriptors of the molecule
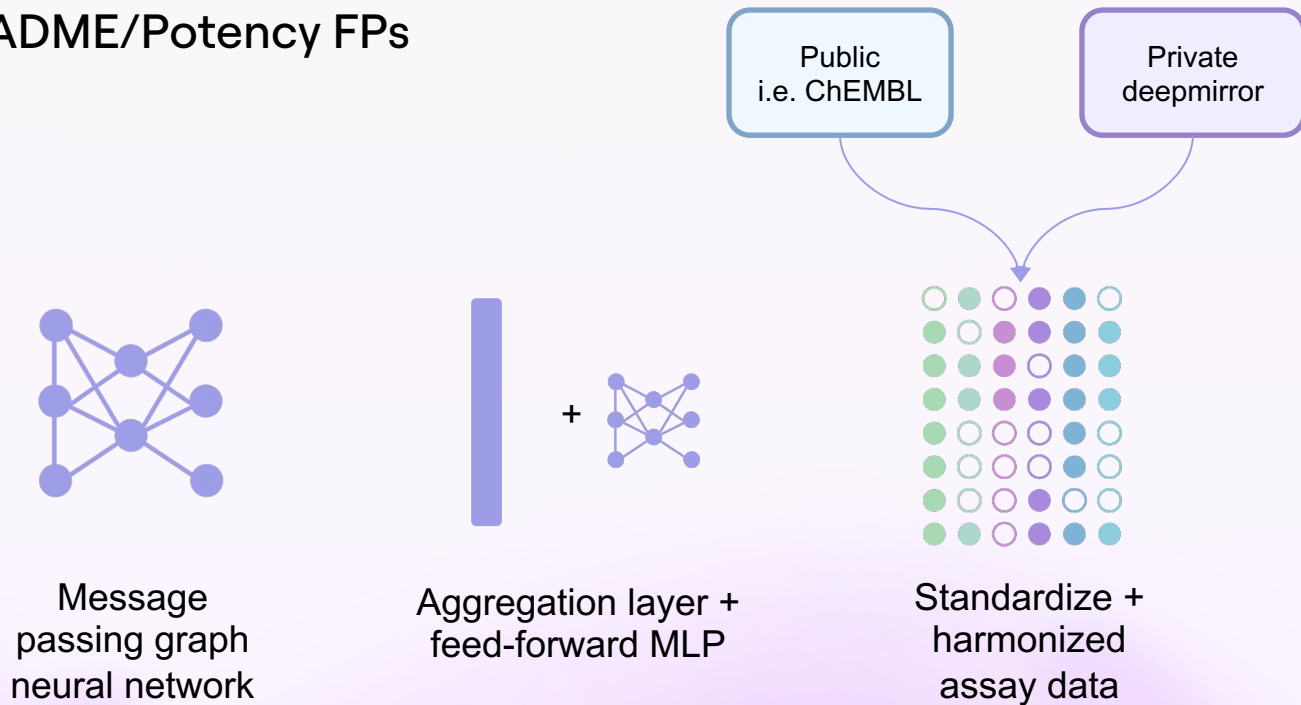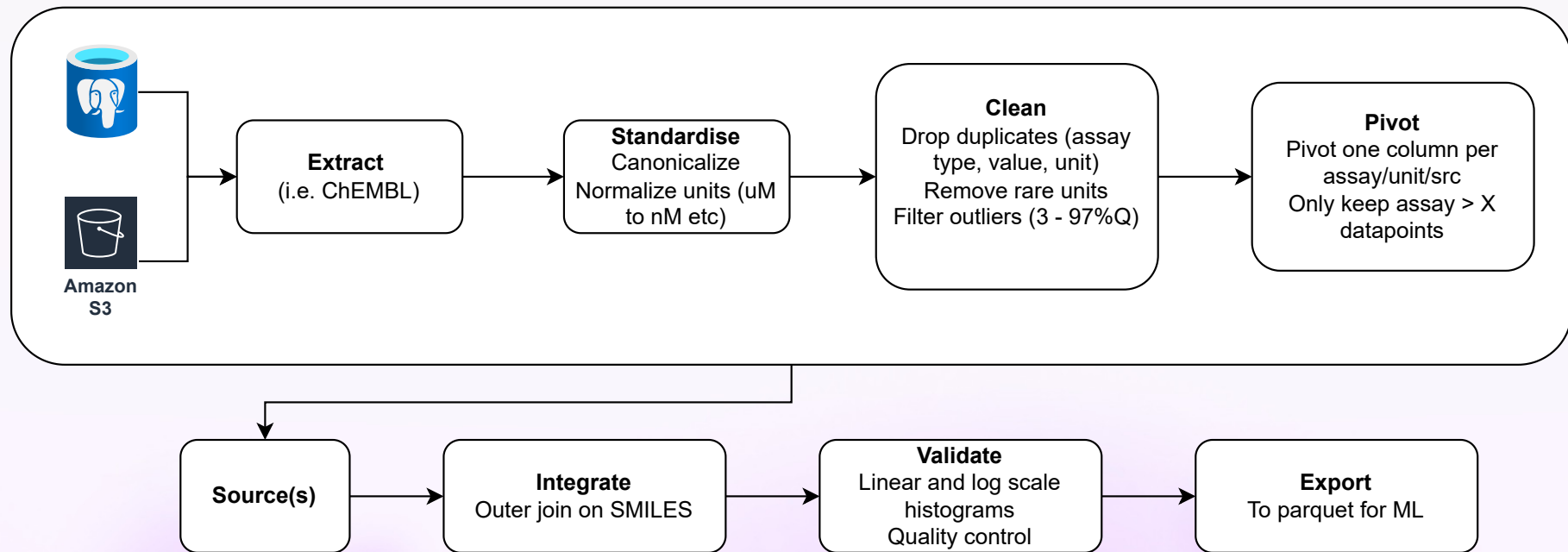
ChEMBL

PubChem

Does training on **assay outcomes (ADMET, potency, selectivity)** make learned fingerprints more predictive for real-world drug discovery?
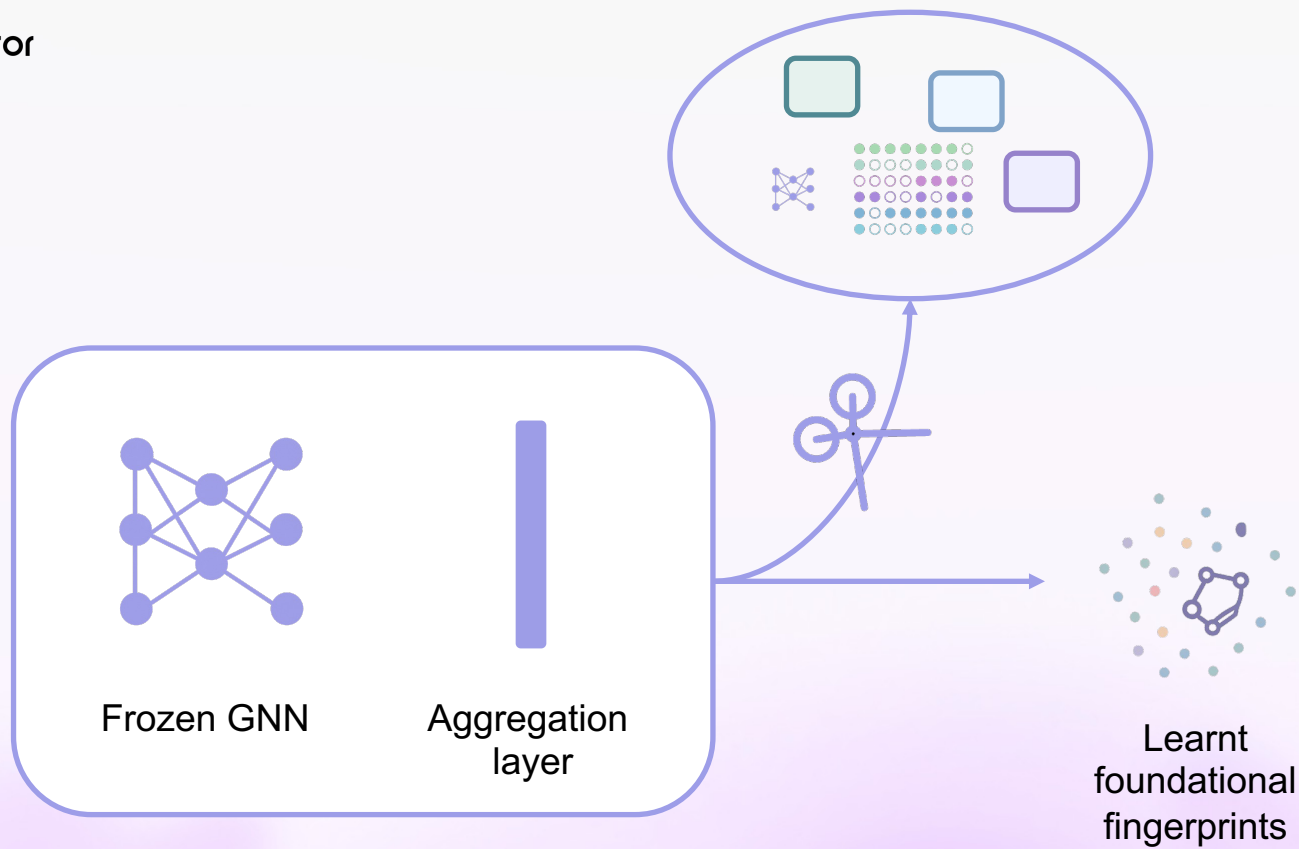
# Combining heterogeneous data

| Molecule | assay-unit-source | assay-unit-source | HLM-hr-labA | HLM-hr-labB | P1-pic50-orgA | P2-pic50-orgB |
|---|---|---|---|---|---|---|
| C=CC(=O… | | 7 | | | | 12 |
| C(=O.. | | | | 8 | 180 | |
| C=CC(=O… | 1.2 | | | | | |
| CCN… | | | 143 | | | |

**deepmirror**

# Pre-training ADME/Potency FPs

Public
i.e. ChEMBL

Private
deepmirror

Message
passing graph
neural network

Aggregation layer +
feed-forward MLP

+

Standardize +
harmonized
assay data

![deepmirror logo]

**Extract**
(i.e. ChEMBL)

**Standardise**
Canonicalize
Normalize units (uM
to nM etc)

**Clean**
Drop duplicates (assay
type, value, unit)
Remove rare units
Filter outliers (3 - 97%Q)

**Pivot**
Pivot one column per
assay/unit/src
Only keep assay > X
datapoints

Amazon
S3

**Source(s)**

**Integrate**
Outer join on SMILES

**Validate**
Linear and log scale
histograms
Quality control

**Export**
To parquet for ML

deepmirror

Frozen GNN          Aggregation
                      layer

Learnt
foundational
fingerprints

![deepmirror]

| Method | $R^2$ (CV) |
|---|---|
| USR | |
| RDKit2D | |
| Pharm3DCATS | |
| Pharm3D | |
| Pharm2D | |
| Morgan FP6 | |
| Morgan FP4 | |
| Mordred | |
| Mol2Vec | |
| MapLight | |
| MTLBERT | |
| MACCS | |
| InfoMax2D | |
| FragECFP | |
| ElectroShape | |
| Desc3D | |
| Chemeleon | |
| ChemBERTa2 | |
| CATS3D | |
| ADMET+Potency FPs | |

- Use ADME and potency data to boost performance on programs.

- Across many programs, the best features depend on the specific context, there isn't a single one that works best for all.

- ADME-based fingerprints are often selected.

- Classical physics-based methods can perform just as well.

deepmirror

Thank you:

Daniel Crusius
+ rest of team

Sign up

Thank you

AMPHISTA
THERAPEUTICS

FERRARIAE UNIVERSITAS
13 · 91
EX LABORE FRUCTUS

HES
PHARMA

…and more

Guillaume Godin &
Jackson Burns

Reach out: ryan@deepmirror.ai