

# Integrating Machine Learning and Quantum Chemistry for Micro-pK<sub>a</sub> Predictions

---

OMRI ABARBANEL

DR. GEOFFREY HUTCHISON

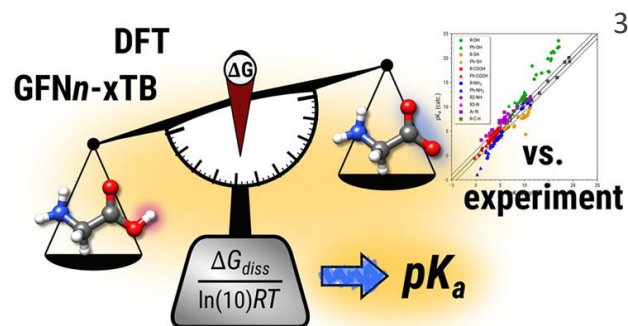
UNIVERSITY OF PITTSBURGH, DEPARTMENT OF CHEMISTRY

ACS FALL MEETING 2023

AUGUST 14, 2023

# Current SOTA

## Quantum Mechanical Models



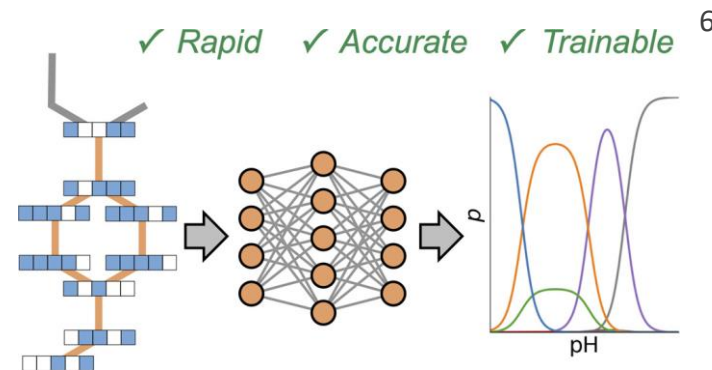
Model	RMSE Range
Jensen 2017 <sup>1</sup>	0.90-1.00
Grimme 2018 <sup>2</sup>	0.68-1.01
Grimme 2021 <sup>3</sup>	0.84-2.68

<sup>1</sup> 10.1021/acs.jpca.6b10990

<sup>2</sup> 10.1007/s10822-018-0145-7

<sup>3</sup> 10.1021/acs.jpca.1c03463

## Machine Learning Models



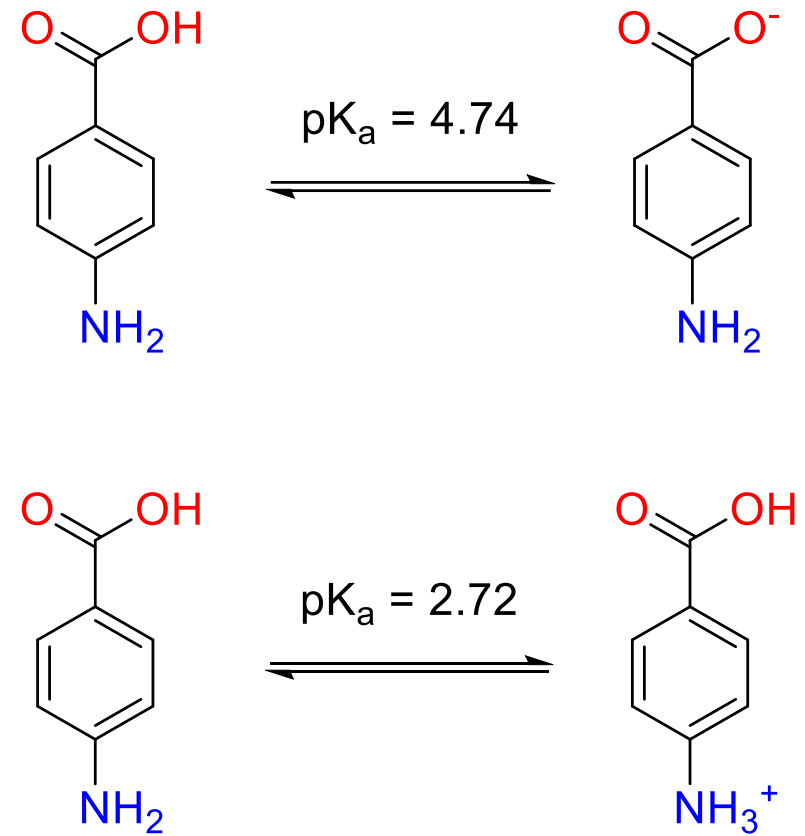
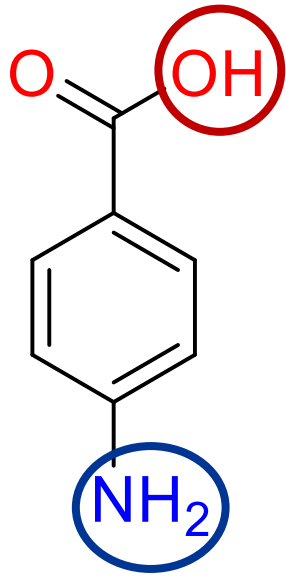
Model	RMSE Range
Czodrowski 2020 <sup>4</sup>	0.79-1.51
Langer 2022 <sup>5</sup>	0.97-1.13
Epik 7 2023 <sup>6</sup>	0.54-1.01

<sup>4</sup> 10.12688/f1000research.22090.2

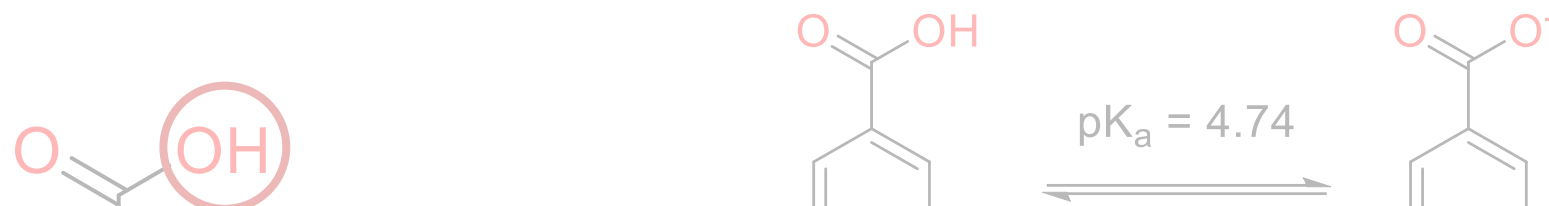
<sup>5</sup> 10.3389/fchem.2022.866585

<sup>6</sup> 10.1021/acs.jctc.3c00044

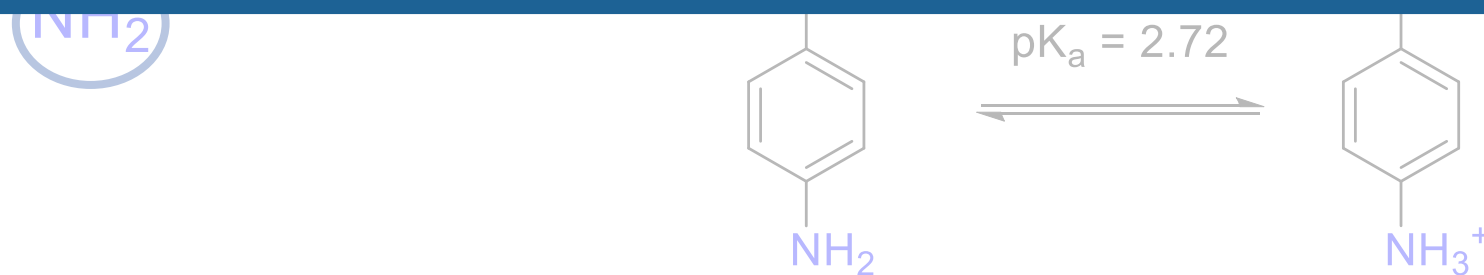
# Micro-pK<sub>a</sub>



# Micro-pK<sub>a</sub>



Can we build a more accurate model to predict micro-pK<sub>a</sub> values by integrating GFN2 features?



# Workflow

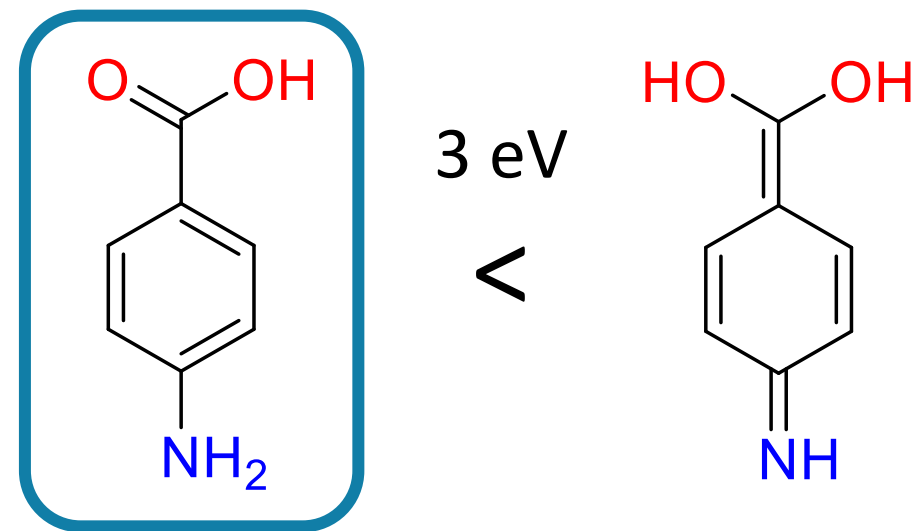


# Workflow



# Tautomer Search

- Tautomer enumeration using RDKit
- GFN2 optimization
  - ALPB implicit solvation model in water
- Lowest energy tautomer is chosen



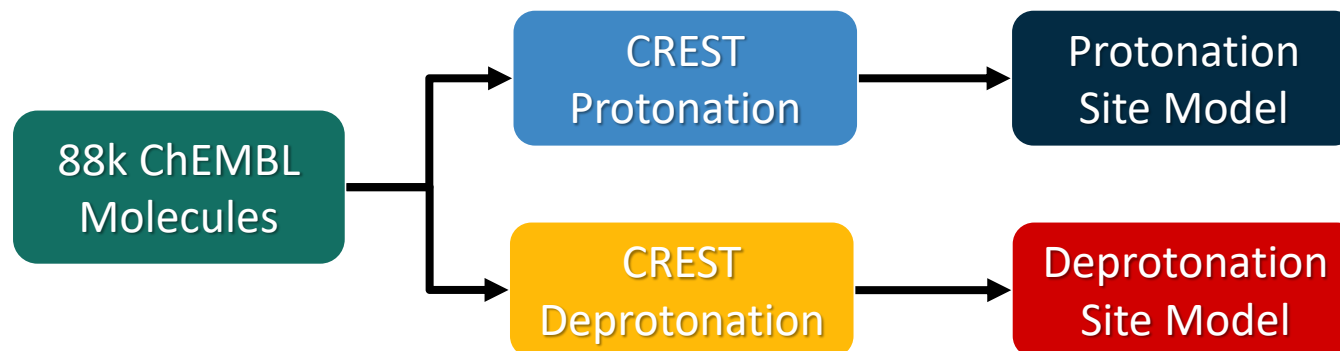
# Workflow





# Reaction Site Enumeration

- CREST<sup>1</sup> protonation \ deprotonation site search<sup>2</sup>
- Two surrogate GNN models
  - Protonation site prediction
  - Deprotonation site prediction



<sup>1</sup> 10.1039/C9CP06869D

<sup>2</sup> 10.1002/jcc.24922

# Graph Representation

## RDKit Features

## GFN2 Features

### Atom

- Atom Type
- No. Heavy Neighbors
- Formal Charge
- Hybridization
- Is in Ring
- Is Aromatic
- Atomic Mass
- VDW Radius
- Cov. Radius
- Chirality
- No. Hydrogens
- Is HBA\HBD

- Partial Charge
- Coord. Number
- Polarizability
- Fukui Indices

### Bond

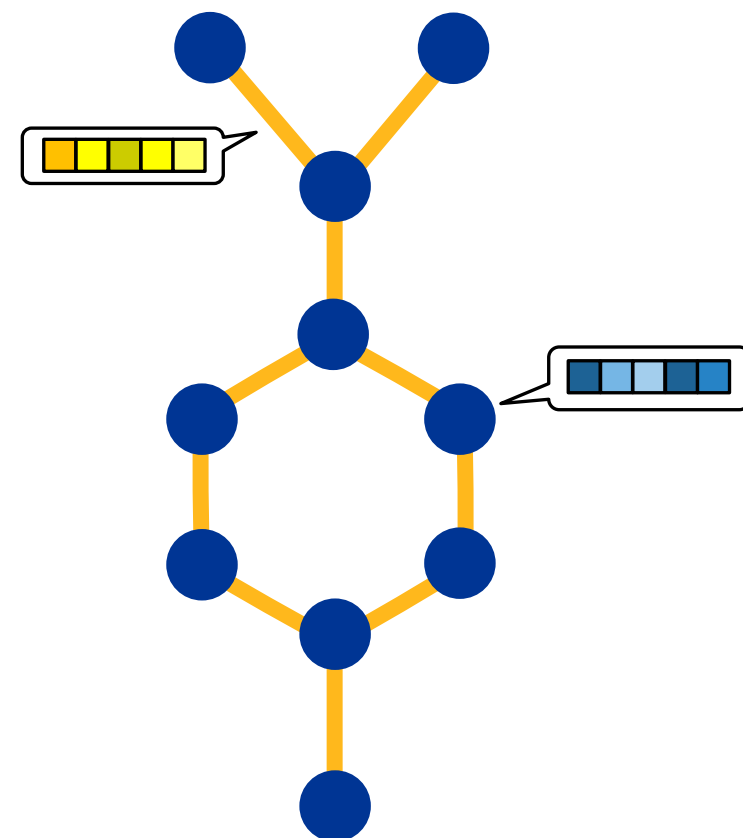
- Bond Type
- Is Conjugated
- Is in Ring
- Stereochemistry

- Bond Order

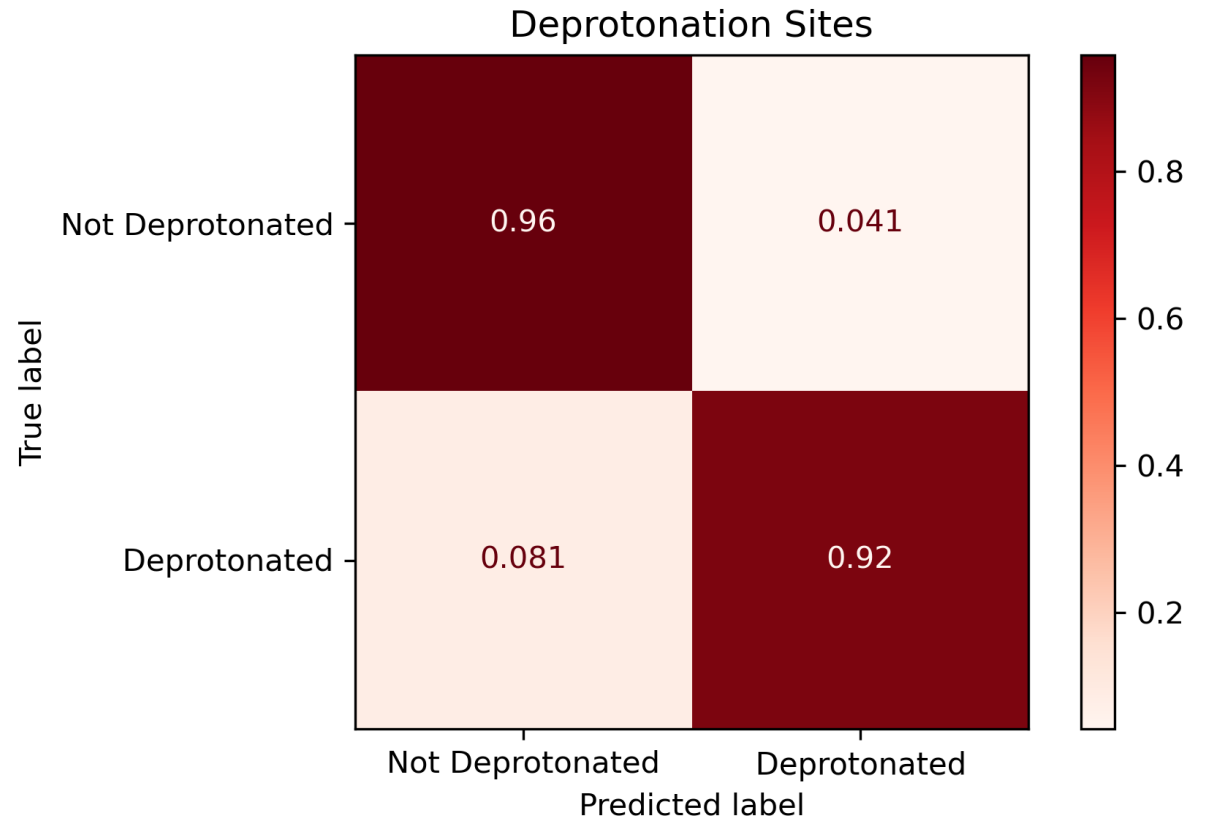
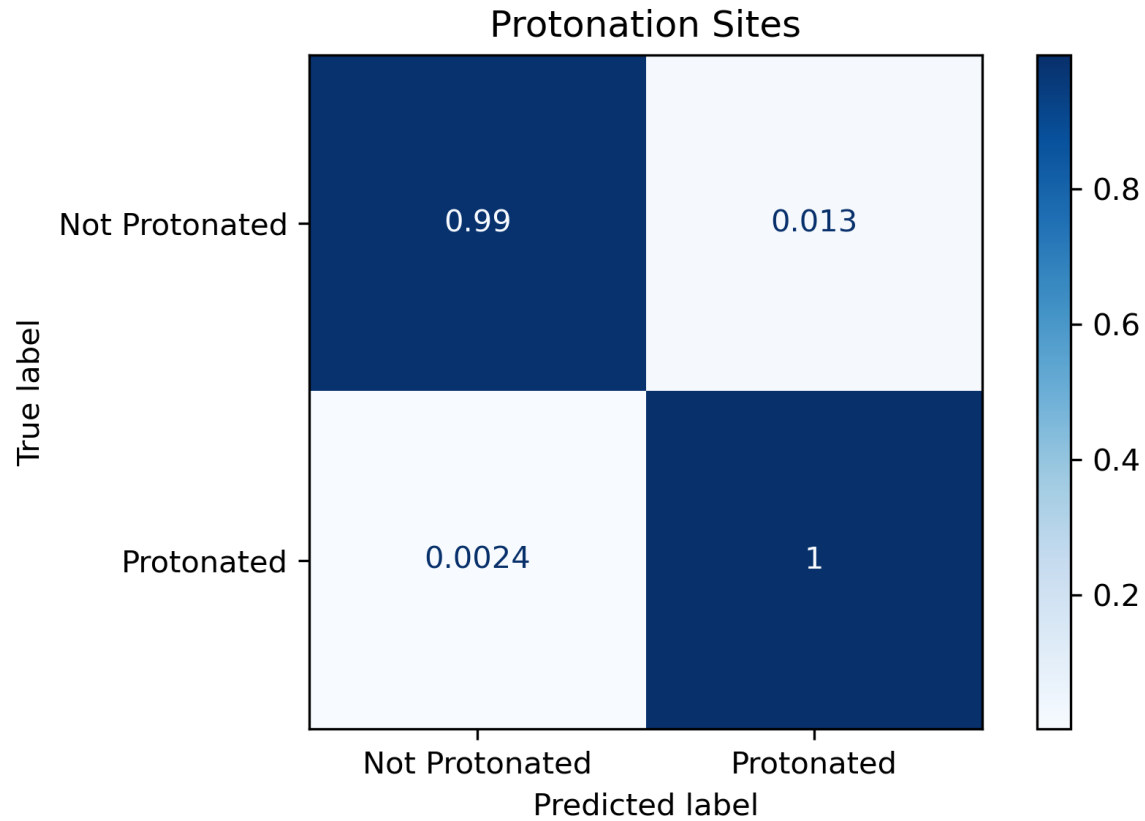
### Molecule

- Radius of Gyration
- Sphericity
- Asphericity
- Eccentricity
- $sp^3$  Fraction

- Charge
- $\Delta E_{\text{ionization}}$



# Site Prediction Models

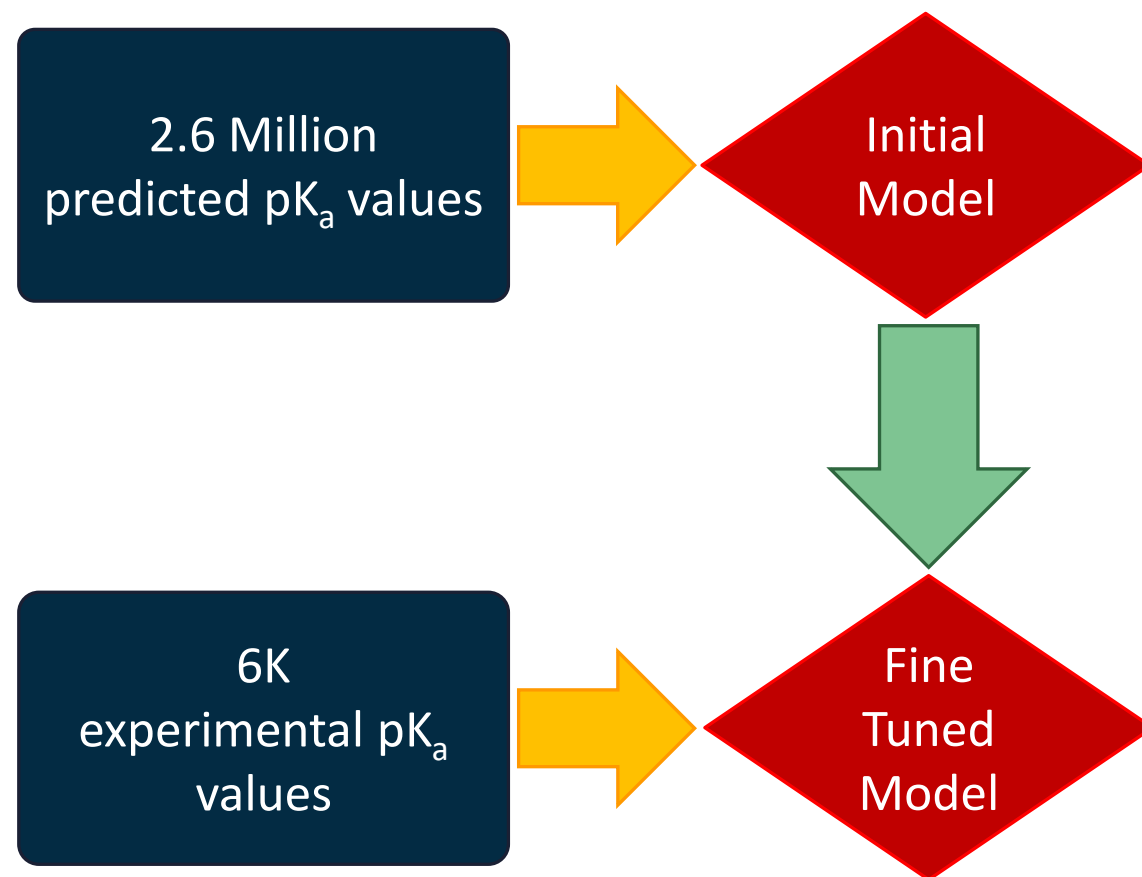


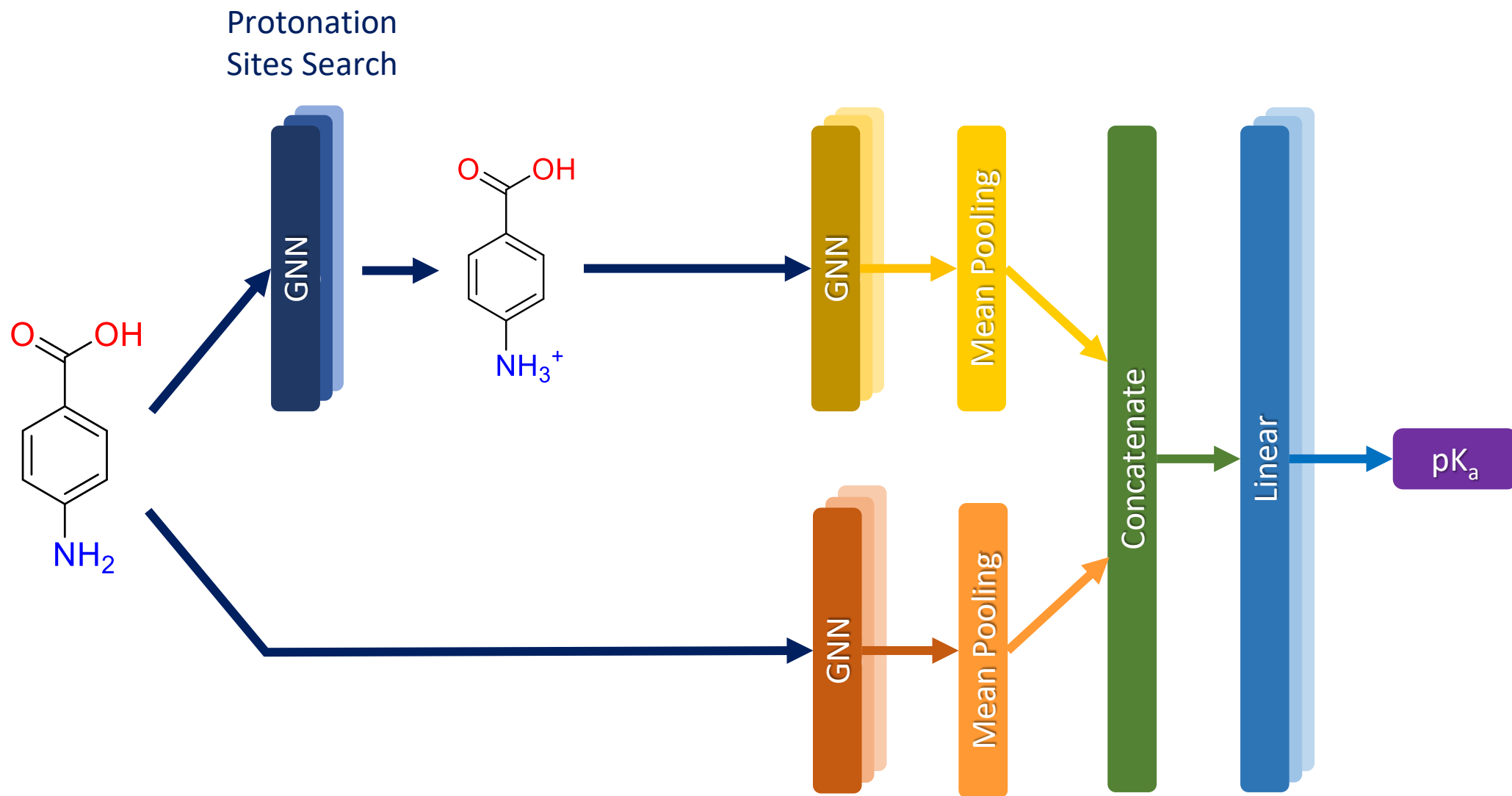
# Workflow

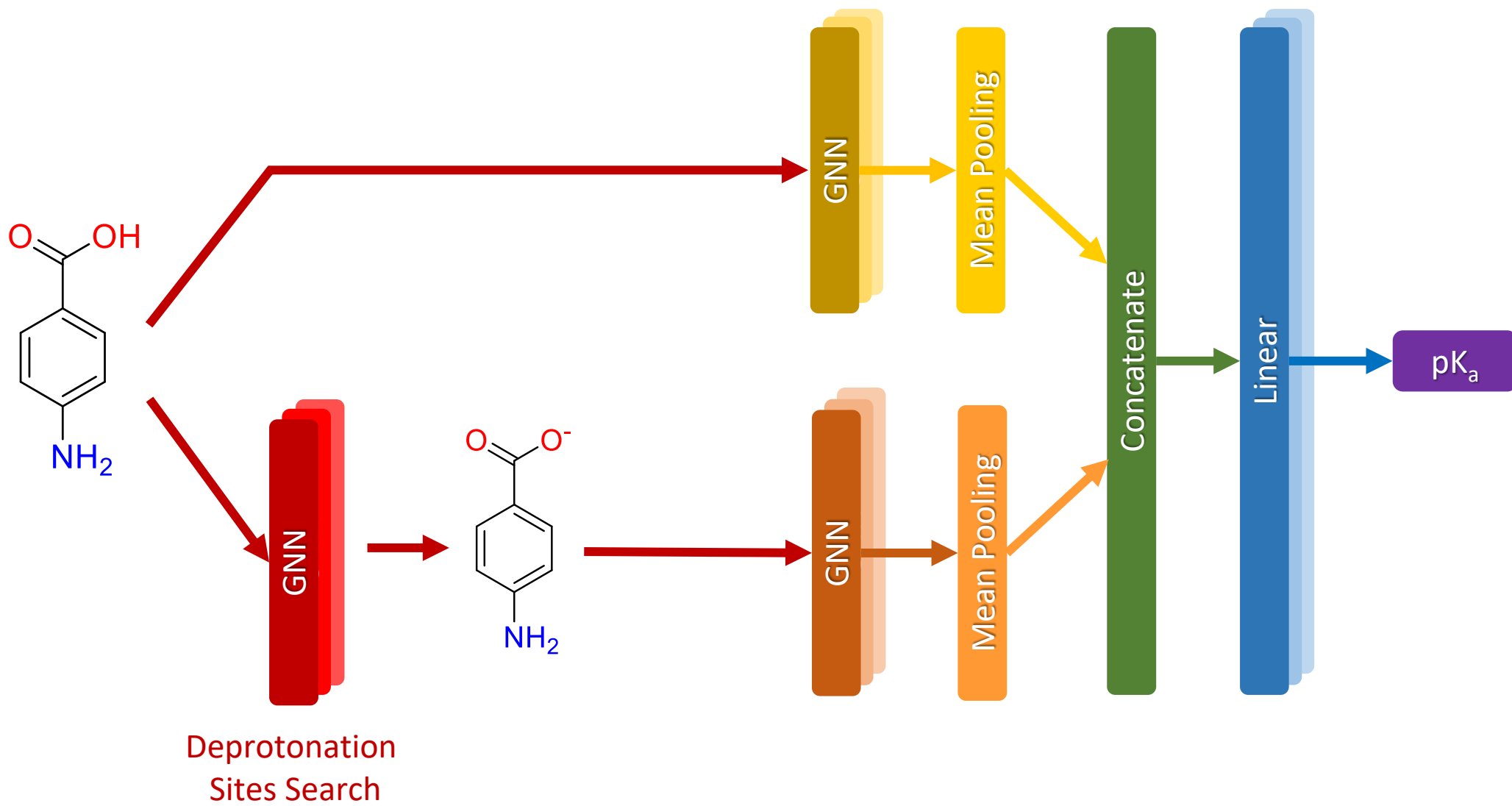


# Micro-pK<sub>a</sub> Prediction Model

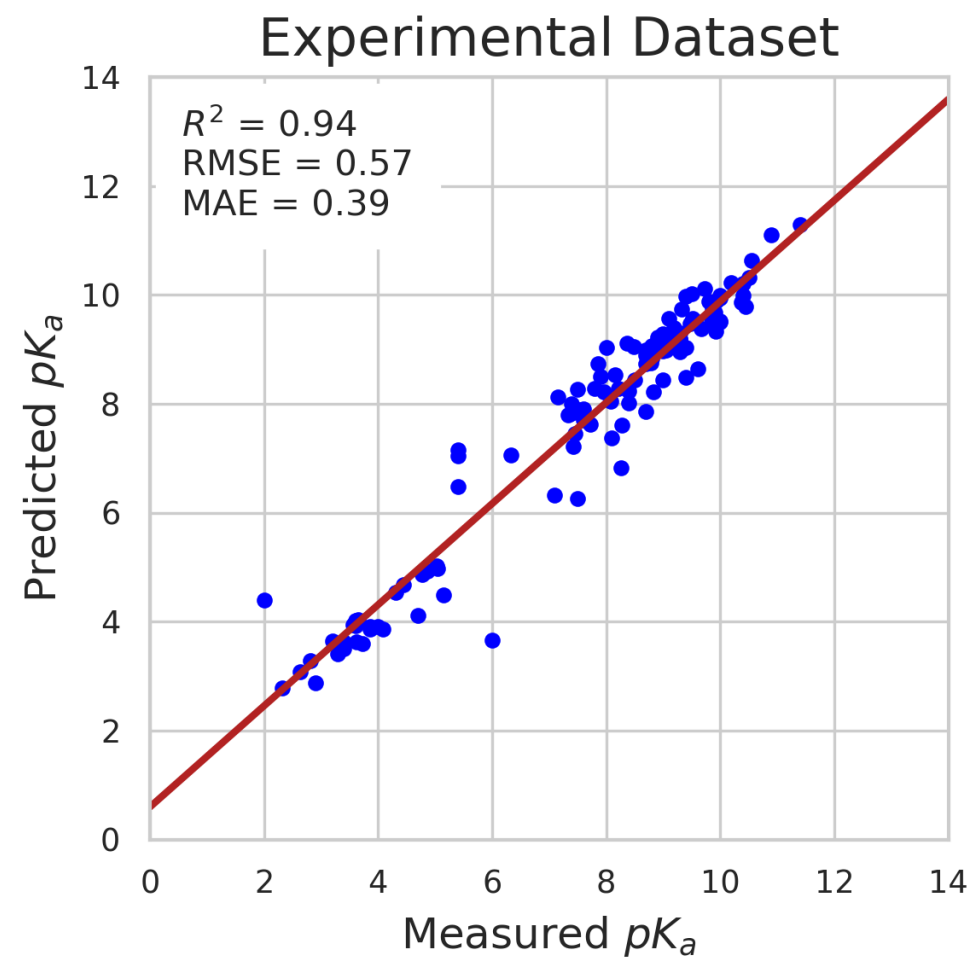
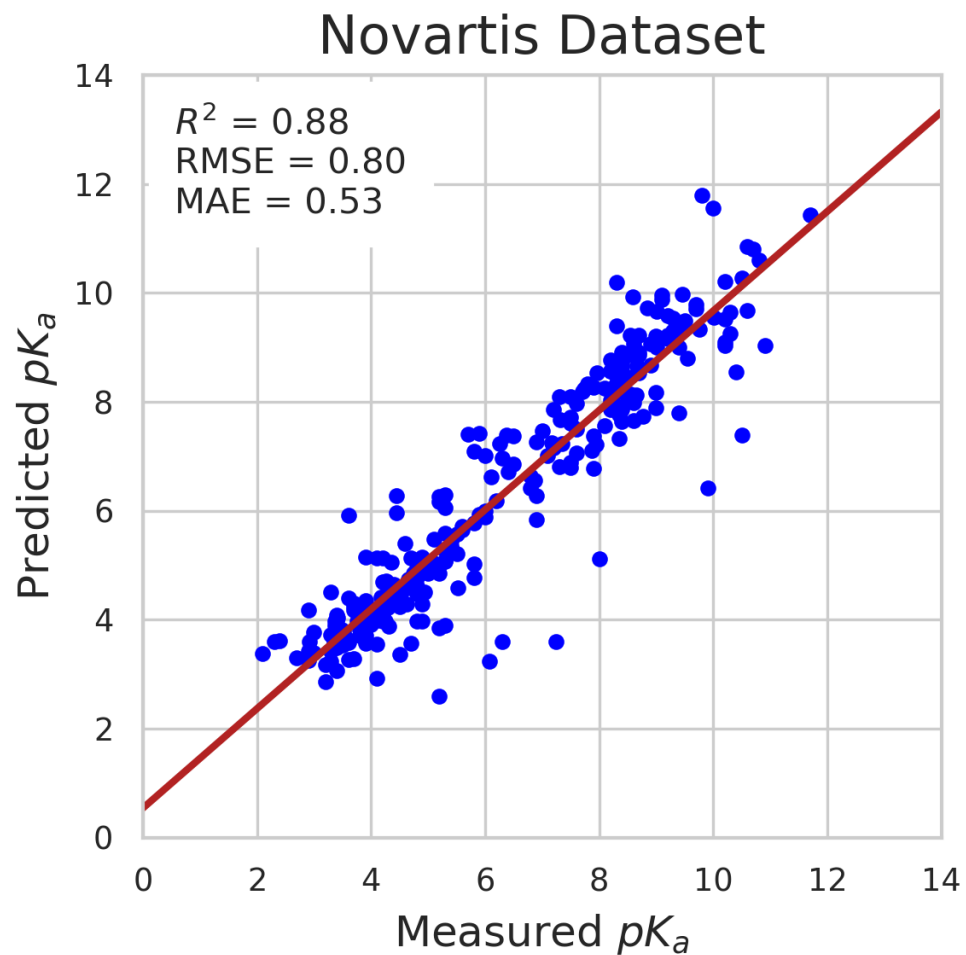
- Problem – Not enough experimental data
  - ~6,000 molecules from the Czodrowski lab
- Solution – Transfer Learning using GNN
  - ~1.5 Million molecules from ChEMBL
  - ~2.6 Million ChemAxon predicted pK<sub>a</sub> values





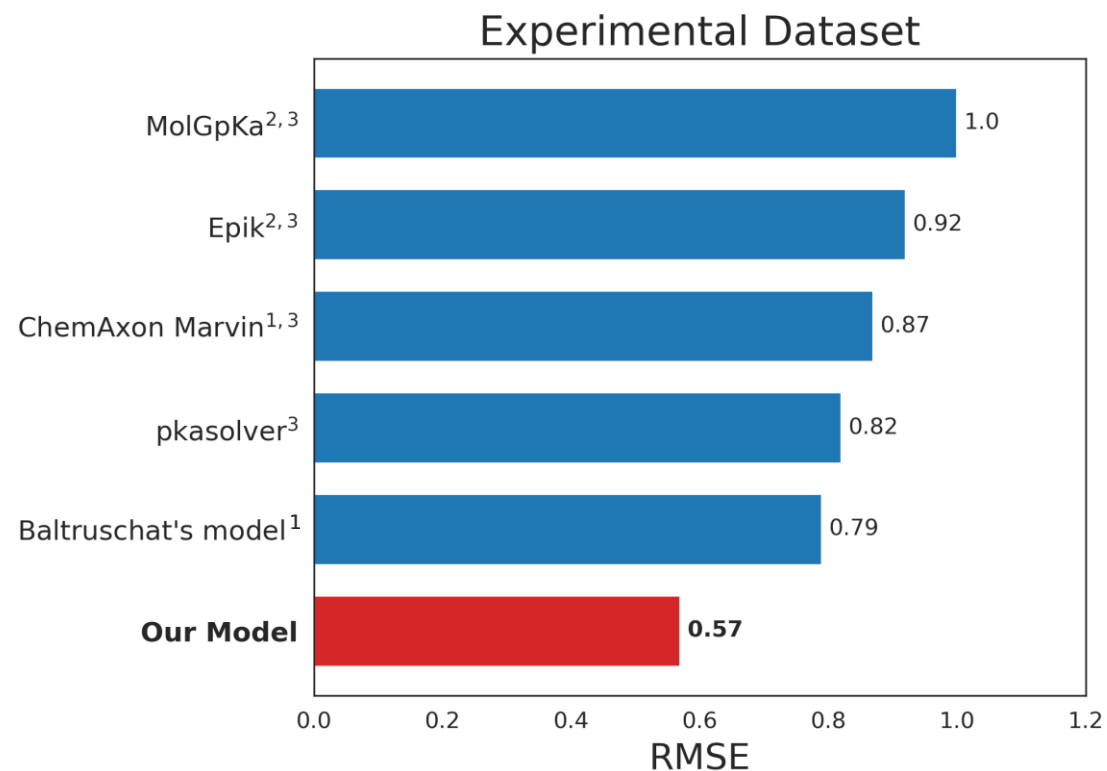
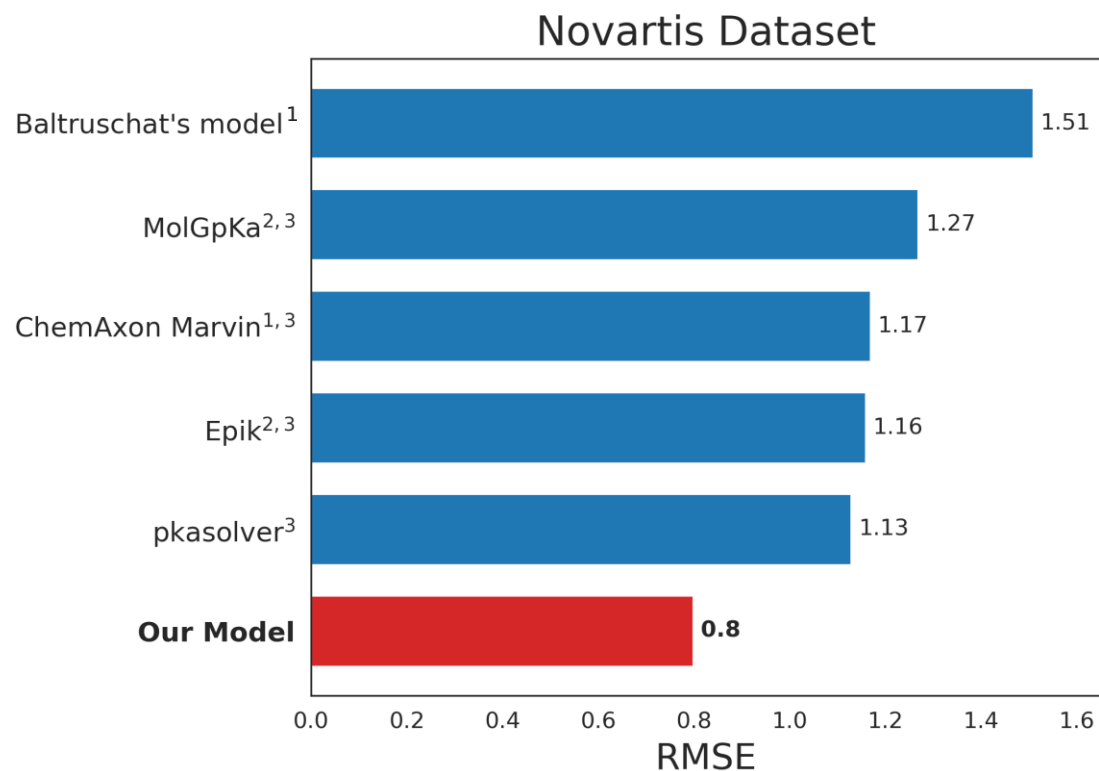


# Test Datasets





# Test Datasets



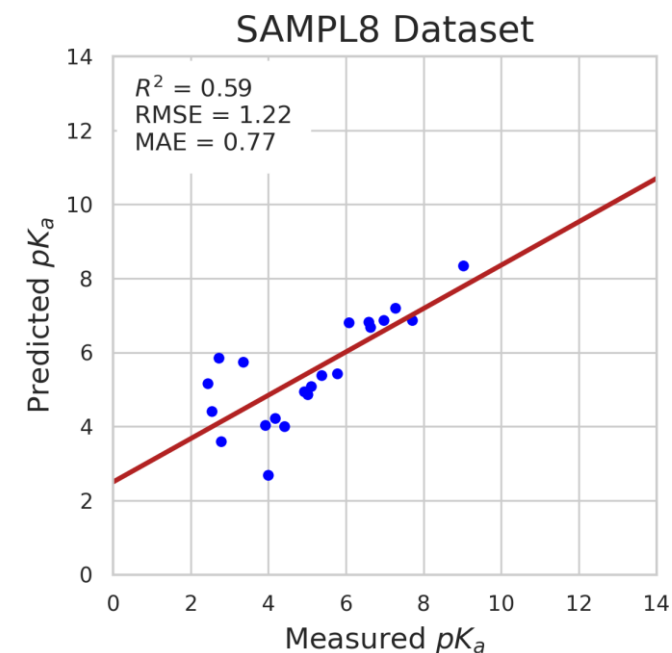
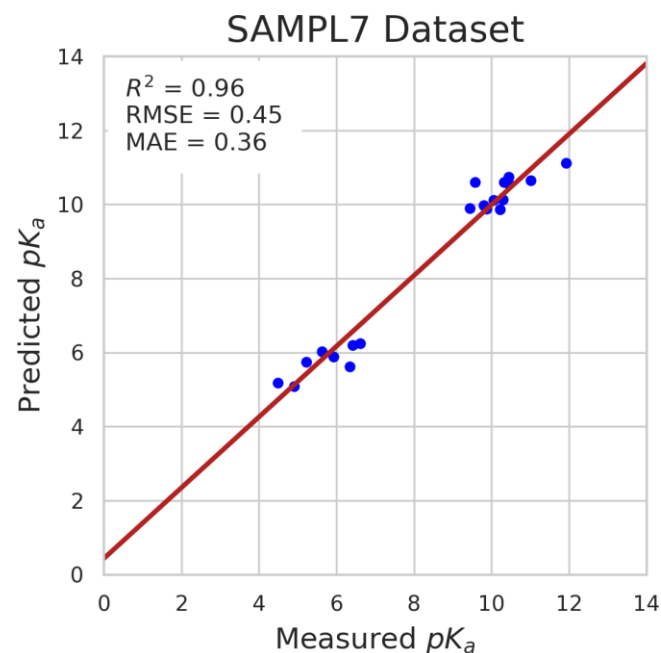
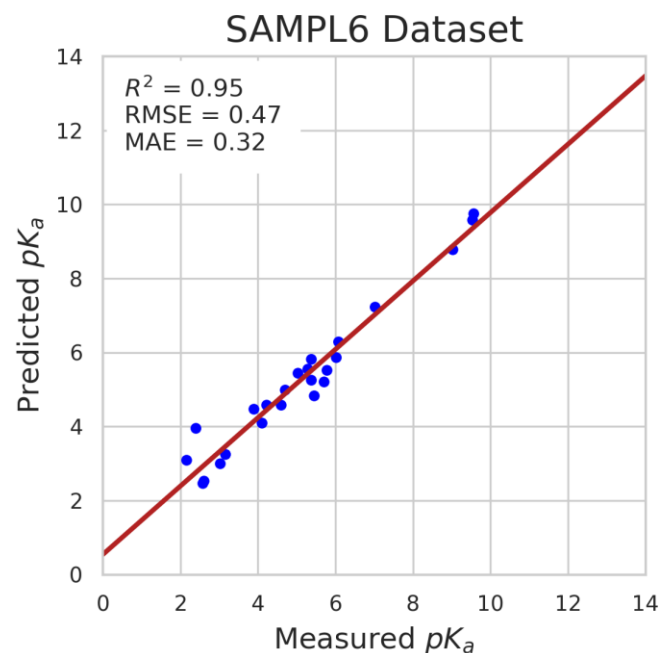
<sup>1</sup> 10.12688/f1000research.22090.2

<sup>2</sup> 10.1021/acs.jcim.1c00075

<sup>3</sup> 10.3389/fchem.2022.866585

<sup>4</sup> 10.1016/j.apsb.2022.11.010

# SAMPL Challenges



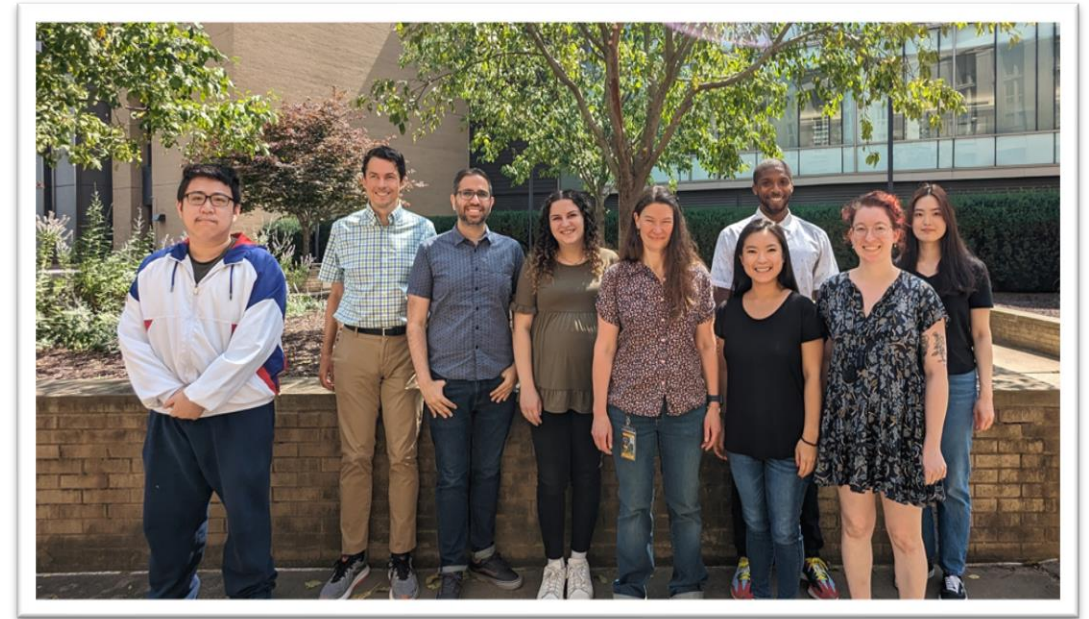
	<b>SAMPL6<sup>1</sup></b>	<b>SAMPL7<sup>2</sup></b>	<b>SAMPL8<sup>3</sup></b>
	RMSE	RMSE	RMSE
Best Submission	0.68	0.71	1.45
<b>Our Model</b>	<b>0.47</b>	<b>0.45</b>	<b>1.22</b>

<sup>1</sup> 10.5281/zenodo.2651393, <sup>2</sup> 10.5281/zenodo.5637494, <sup>3</sup> 10.5281/zenodo.7535037

# Conclusions

---

- First free, open model with micro-pK<sub>a</sub> predictions with RMSE below 0.8
- Adding semi-empirical QM features can help with micro-pK<sub>a</sub> predictions
  - Comparisons with ML models (e.g., ANI / AIMNET) to come
- More experimental data is needed (please!)
  - Transfer learning is a viable way to remedy that
- ChemRxiv manuscript, code and data to come soon



# Thank You

QUESTIONS?



**PittCRC**  
Center for Research Computing

