# Schrödinger
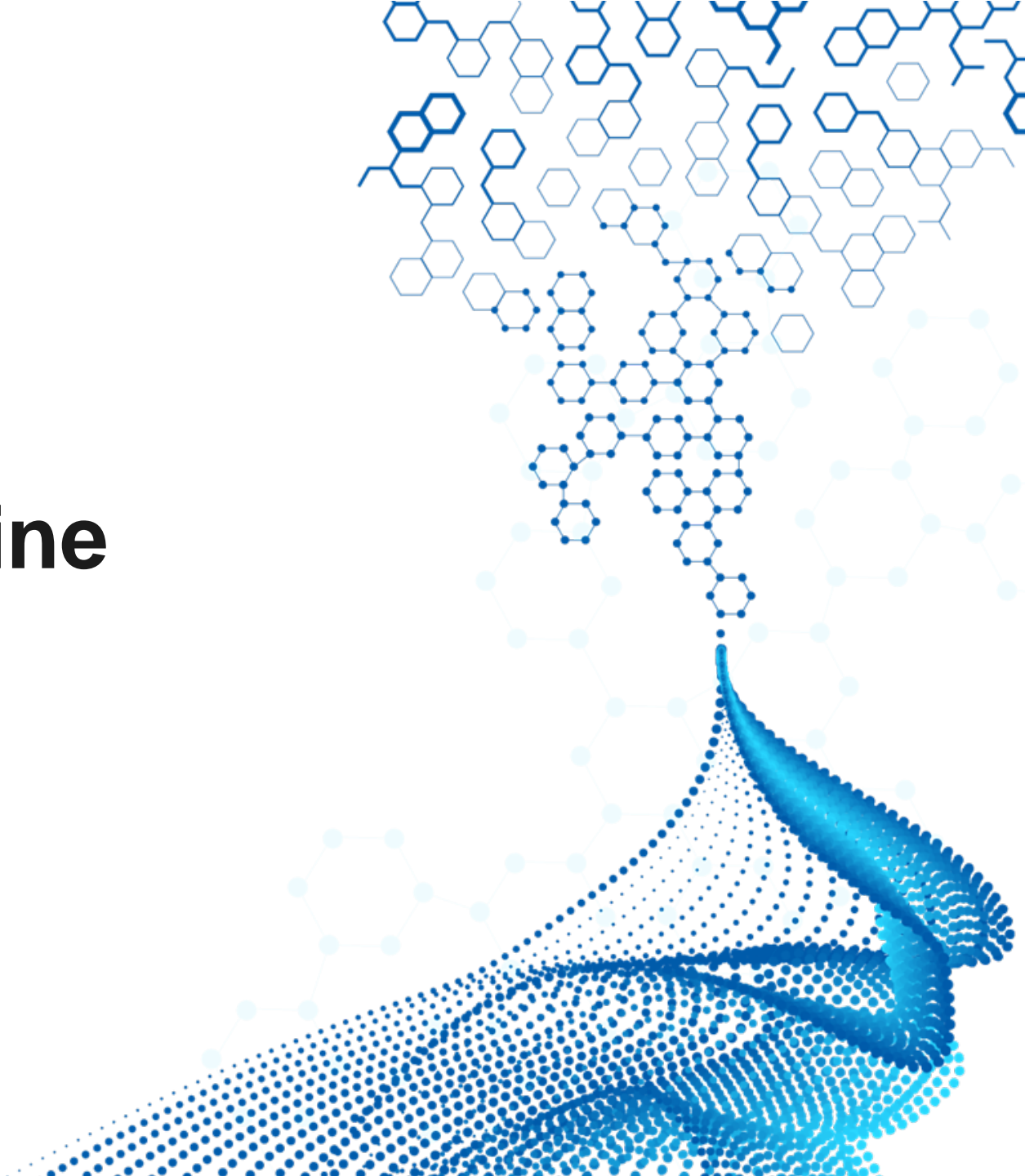
# Data-driven materials innovation: where machine learning meets physics
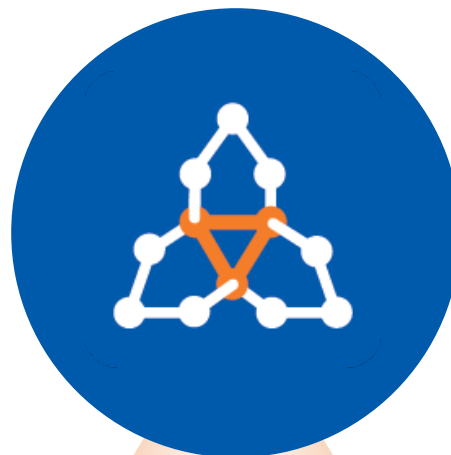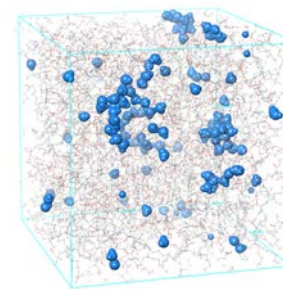
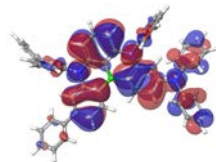Anand Chandra
Product Manager, Materials Science Informatics
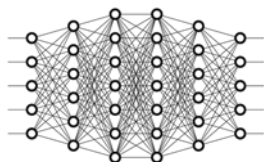chandras@schrodinger.com
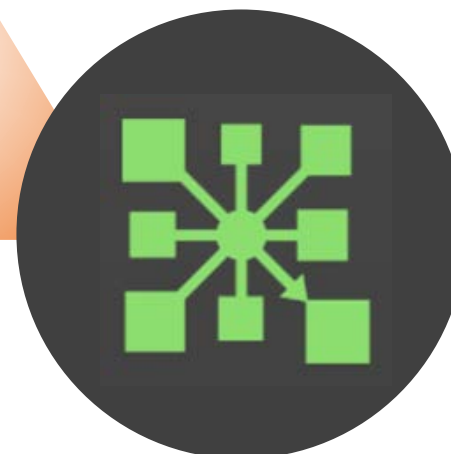
October 31st 2023

# Machine Learning for Materials Design/Discovery at Schrödinger
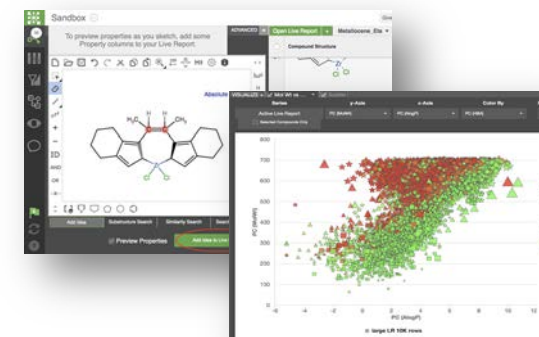
Expertise in physics-based simulation
and domain knowledge

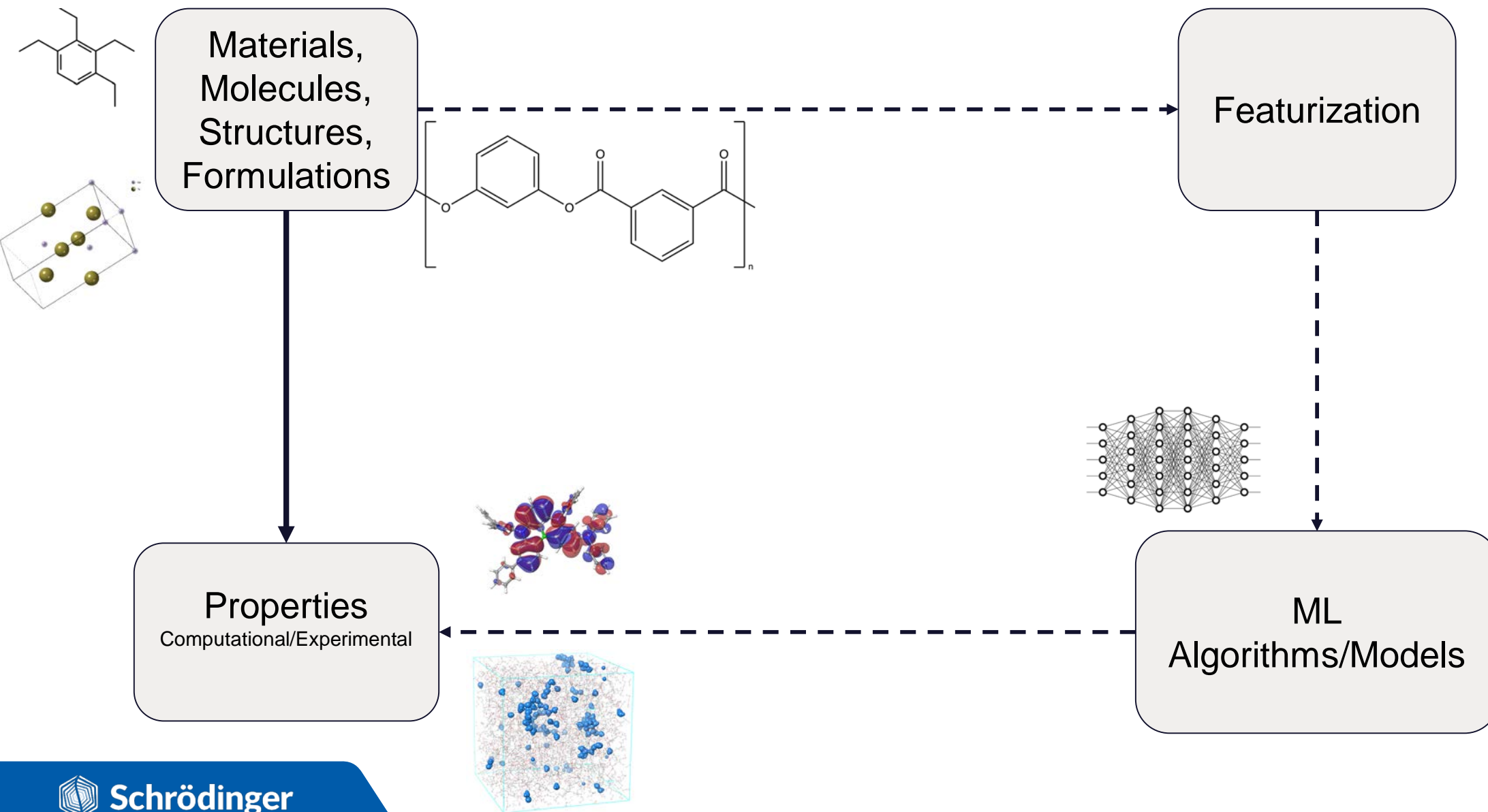Latest machine learning technology for
materials chemistry

Enterprise solution for data management
and collaborative ideation

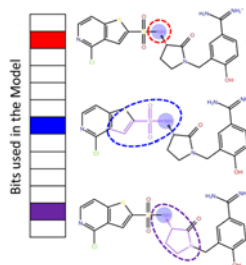# Supervised Learning in Materials Science

# Featurization in Diverse Materials Systems

- Properly featurizing various chemical systems is key to building predictive machine learning models
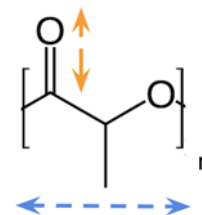
- Small Molecules
  - Physiochemical, topographical descriptors
  - Binary fingerprints (RDKit, Canvas)
  - Graph-based convolution neural networks

- Polymers
  - Taking into account connections between repeat units
  - RDKit fingerprints + customized descriptors

- Periodic Inorganic Solids
  - Element
  - Lattice structure
  - Oxidation state
  - Intercalation descriptors
  - 3D SOAP (with PCA)

- Formulations and Mixtures
  - Composition
  - Chemistry of the components
  - Experimental/Processing conditions

*Schrödinger's Physics-based Simulation Provides Additional Power to Machine Learning*

**QM (Jaguar)**
**Catalysis (AutoRW)**

**MD (Desmond)**
**MD + QM**

**Periodic QM (QE)**

**MD (Desmond)**

**Composition**

GeTe

**Site and Structure**

# Automated Machine Learning and Visualization in Molecular Systems

- Supervised learning with 400+ built-in descriptors
- Integrated as automated HPC-supported workflow



Schrödinger's automated model-building algorithm (AutoQSAR)

Machine Learning with Model visualization



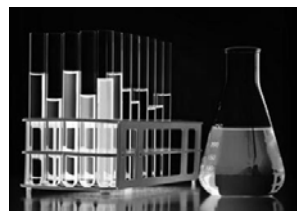Kernel-Based Partial Least Squares: Application to Fingerprint-Based QSAR with Model Visualization

Automated cross-validation for model scoring and ranking

# AutoQSAR for Ionic Liquids

- 392 ionic liquids from the NIST IL Thermo database
- Target Property → Electrical conductivity

# DeepAutoQSAR: Automated Model Selection & Parameter Optimization



Data splits, featurization

Model architecture, descriptor and hyperparameter combinations explored and optimized via Bayesian Optimization

Scoring

Ranking

**Models Sampled**
- Dense Neural Network
- Random Forest Regressor
- XGBoost
- TorchGraphConv
- GCN
- GraphSAGE
- GIN
- TopK
- SAGPool
- EdgePool
- GlobalAttention
- Set2Set
- SortPool

**Consensus Model**
- Prediction = an average of the predictions for 5 best models
- Uncertainty = SD across the 5 predictions

Schrödinger

# Case Study - Redox Flow Batteries

- Design: Homobenzylic ethers (HBE) with oxidation potential in a pre-specified range.

- Oxidation potential of 1,400 HBEs calculated as the initial (training) dataset for machine learning



Figure 2 Graphical illustration of the chemical space of 1,400 HBEs and their computed oxidation potentials ($E^{ox}$). PC-1 and PC-2 represent principle component 1 and 2, respectively.

[1] Doan, Hieu A., Garvit Agarwal, Hai Qian, Michael J. Counihan, Joaquín Rodríguez-López, Jeffrey S. Moore, and Rajeev S. Assary. "**Quantum Chemistry-Informed Active Learning to Accelerate the Design and Discovery of Sustainable Energy Storage Materials**." *Chemistry of Materials* (2020).

Schrödinger

# AutoQSAR vs DeepAutoQSAR Results



Traditional AutoQSAR
Test $R^2$ = 0.94

DeepAutoQSAR
Test $R^2$ = 0.98

- ML models were created for oxidation potential of 1,400 homobenzylic ethers for Redox Flow

- Both AutoQSAR and DeepAutoQSAR offer solid predictive capability.

- The deep-NN-based model (by DeepAutoQSAR) outperforms descriptor-based models for larger (>1000) training set.

Schrödinger

# Chemical Featurization using Physics

- 100+ additional physics-based descriptors by QM-bound properties, repeat-unit chemistry, and crystallinity

- Direct link to AutoQSAR and other workflows within the platform



Molecular descriptors

Periodic descriptors

Polymer descriptors

# Customized Polymer Descriptors Outperform Simple Monomers

## Polymer Descriptors

- Topological torsion fingerprints
- Number of rotatable bonds
- Number of ring atoms
- Fused ring atoms



Tg ML model AutoQSAR monomer

R2=0.898
RMSE=37.939

RMSE = 38 K

Data set of 315 polymers with Tg values

RMSE = 11 K

Train set ($R^2$=0.983)
Test set ($R^2$=0.950)

# Viscosity Dataset for Machine Learning Models

Literature extraction of viscosity
~5,356 viscosities

Single, organic structures

methyl acetate ✓    silver nitrite ✗

Remove extreme $\mu$ and $T$

$\mu$    $T$

Remove positive deviations of $\mu$ vs. $T$

$\mu$    $T$

**4,400 viscosities**

## Distribution of viscosity and temperature

Frequency (log $\mu$): 1400, 1200, 1000, 800, 600, 400, 200, 0
log $\mu$ axis: -1.5, -1.0, -0.5, 0.0, 0.5, 1.0, 1.5

Frequency (Temperature): 1400, 1200, 1000, 800, 600, 400, 200, 0
Temperature (K) axis: 200, 250, 300, 350, 400, 450

## Dataset summary:

- **1,005** unique molecules

- Atomic elements of {H, C, N, O, F, Si, P, S, Cl, Br, and I}

- Viscosity is between **0.10** to **26.52 cP**

- Temperature is between **227 K** to **404 K**

# Quantitative Structure-Property Relationships (QSPR)

**Descriptor-Based**



**Graph neural networks (GNN)-Based**

# Impact of MD-Derived Simulation Descriptors

MD Snapshot at T = 298 K

Methyl acetate

~1-2 hrs

**Eight MD Descriptors**
- Heat of vaporization
- Density
- Hansen solubility parameters
- Root-mean-square displacement

## Performance with MD Descriptors



## Learning curve with and without MD Descriptors



## Main takeaways:

- Inclusion of MD Descriptors lowers test set RMSEs

- MD descriptors are most useful at small training size (<1,000)

*Scheduled for 23-4*

# Impact of MD-Derived Simulation Descriptors

MD Snapshot at T = 298 K

Methyl acetate

~1-2 hrs

**Eight MD Descriptors**
- Heat of vaporization
- Density
- Hansen solubility parameters
- Root-mean-square displacement

## Performance with MD Descriptors



## Learning curve with and without MD Descriptors



**Main takeaways:**

- Inclusion of MD Descriptors lowers test set RMSEs

- MD descriptors are most useful at small training size (<1,000)



LGBM (2D and MD)

**Which descriptors were most useful for viscosity?**

15

# Machine Learning Optoelectronics Properties with DFT descriptors

# Database of Optical Properties of Organic Compounds

- Experimental dataset of **20,236** combinations of **7,016** chromophores in **365** solvents



Joung, Joonyoung F., et al. "Deep learning optical spectroscopy based on experimental database: potential applications to molecular design." *JACS Au* 1.4 (2021): 427-438.

# Benchmark of DFT Descriptors

- Combining 2D and DFT descriptors leads to state-of-the art performance



*Scheduled for 23-4*

*The technical features and projected timeline presented on this slide is for discussion purposes only. Such planned or potential capabilities are subject to change at any time.*

**MatSci-ML Model: Neural Network**

Features: 2D Descriptors of Chromophore (+DFT Features) + Dielectric constant of Solvent

# Feature Importance Analysis



λmax Absorption — bar chart (Mean |SHAP|):
- r_matsci_optelec_Blue_Area
- r_matsci_optelec_S1_at_S0_(eV)
- rdkit_descr_MaxAbsPartialCharge_c
- i_matsci_optelec_Lmax_(nm)
- r_matsci_optelec_Scaled_Gap_(eV)

λmax Emission — bar chart (Mean |SHAP|):
- r_matsci_optelec_Red_Area
- r_matsci_optelec_S1_at_S0_Transition_Dipole_(D)
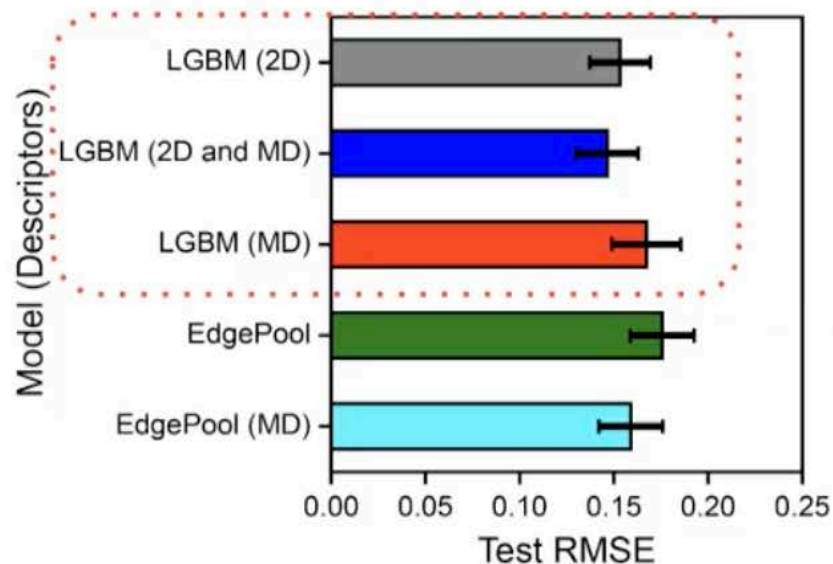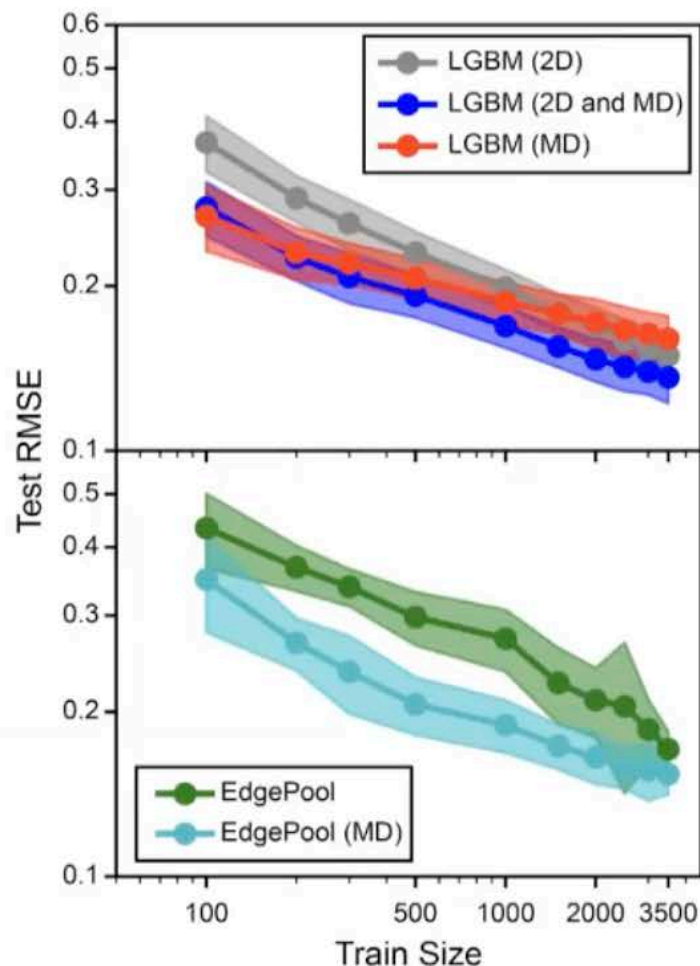- Dielectric
- r_matsci_optelec_S1_at_S0_(eV)
- r_matsci_optelec_Scaled_Gap_(eV)

## DFT Descriptors

- **i_matsci_optelec_Lmax_(nm) \***
- r_j_Final_Energy
- r_j_Gas_Phase_Energy
- r_j_HOMO \*
- r_j_LUMO \*
- r_j_QM_Dipole_(debye)
- r_matsci_optelec_Blue_Area
- r_matsci_optelec_Dipole_(D)
- r_matsci_optelec_Green_Area
- r_matsci_optelec_Oxidation_Potential_(eV)
- r_matsci_optelec_Red_Area
- r_matsci_optelec_Reduction_Potential_(eV)
- **r_matsci_optelec_S1_at_S0_(eV) \***
- **r_matsci_optelec_S1_at_S0_Transition_Dipole_(D) \***
- r_matsci_optelec_S2_at_S0_(eV)
- r_matsci_optelec_S2_at_S0_Transition_Dipole_(D)
- r_matsci_optelec_S3_at_S0_(eV)
- r_matsci_optelec_S3_at_S0_Transition_Dipole_(D)
- **r_matsci_optelec_Scaled_Gap_(eV) \***
- r_matsci_optelec_Scaled_HOMO_(eV)
- r_matsci_optelec_Scaled_LUMO_(eV)

**\* most impactful via feature importance**

# Machine Learning for Volatility of Organic Molecules

# Evaporation/Sublimation of Organic Molecules

- Training data:
  - 1,184 organic molecules containing C, O, Cl, N, Si, Br, S, F, P, I, B, As, Se
  - **12,169** experimental $(p, T)$ datapoints
  - Pressure ranges from 1 Torr to 30 atm

- Generate 200 chemical descriptors and 1000 Morgan Fingerprints for each molecule from its 2D sketch
  - Examples of descriptors: molecular weight, solvent-accessible volume, max partial charge on atoms, electrotopological state descriptors …

- Log($p$) was used as an additional descriptor and ML model was trained to predict 1/$T$

# Benchmarking ML Algorithms

Top-performing machine learning algorithms:

- Light Gradient Boosting Machine (LightGBM)
  - RMS error ±8°C
- Multi-Layer Perceptron (neural network)
  - RMS error ±2°C

Most literature QSPR models for boiling points of diverse organic molecules have errors ±18°C

"Quantitative structure-property relationships for prediction of boiling point, vapor pressure, and melting point", J. C. Dearden, Environmental Toxicology and Chemistry, 22, 1696–1709 (2003).

Best neural network gives RMS-error ±5°C and mean absolute error ±4°C

"Boiling point and critical temperature of a heterogeneous data set: QSAR with atom type electrotopological state indices using artificial neural networks", L. H. Hall & C. T. Story, J. Chem. Inf. Comput. Sci. 36, 1004–1014 (1996).



Model score 0.997

Train set (10952 points, $R^2$=0.999, $RMSE$=0.000)
Test set (1217 points, $R^2$=0.998, $RMSE$=0.000)

# Prediction of Pressure-Temperature Relationships

Performance of model on sample molecules *outside* training set

# Applications of Volatility Machine Learning

- Atomic Layer Deposition / Chemical Vapor Deposition

- Thermal evaporation & jet-printing (Organic LED)

- Flavors & fragrances

- Equation of state for petroleum fluids

- Refrigerants

- Membrane separation/distillation

- Volatile Organic Compound Pollutants

- Explosion hazards

# Machine Learning for Inorganic 3D Crystal Structures

# Transparent Conducting Oxide Band Gap ML

- 3000 periodic structures containing indium, aluminum, gallium and oxygen

- These materials have applications in display devices and solar-cells

- Dataset was obtained from [NOMAD 2018 Kaggle challenge](#) on creating ML models for properties of transparent conducting oxides

- ML models for **Band Gap** were created using DeepAutoQSAR

- Composition (matminer) and 3D SOAP descriptors were used

# DeepAutoQSAR Results



R² on test set = 0.886

With 139 matminer descriptors

**AutoQSAR/DeepChem**

**Choose task:** ○ Build model  ● Make Predictions

Model file: r/chandras/Schrodinger/NOMAD2018.prj/deepchem_no_SOAP/deepchem_no_SOAP_model.qzip  [ Browse... ]

**Model Summary**

| Mode: | Regression |
|---|---|
| kendall_tau: | 0.7967 |
| mae: | 0.2405 |
| median_ae: | 0.1598 |
| r2: | 0.8869 |
| rmse: | 0.3542 |
| time: | 2021-05-07T21:56:03 |

[ View Full Report ]

R² on test set = 0.939

With 139 matminer descriptors + 10 SOAP-PCA descriptors

**AutoQSAR/DeepChem**

**Choose task:** ○ Build model  ● Make Predictions

Model file: andras/Schrodinger/NOMAD2018.prj/deepchem_all_descr/deepchem_all_descr_model.qzip  [ Browse... ]

**Model Summary**

| Mode: | Regression |
|---|---|
| kendall_tau: | 0.8632 |
| mae: | 0.1480 |
| median_ae: | 0.0750 |
| r2: | 0.9392 |
| rmse: | 0.2597 |
| time: | 2021-05-07T21:41:31 |

[ View Full Report ]

Schrödinger

# Machine Learning Property Prediction Panel

- The following properties/models are currently available

    - Volatility of organic molecules (Both Boiling Point and Vapor pressure)
    - Volatility of organometallic molecules
    - Polymer Tg
    - Frequency dependent Df
    - Frequency dependent Dk
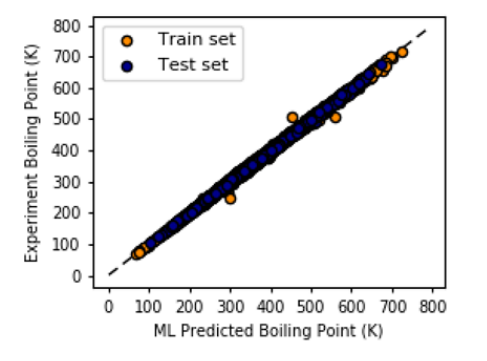    - Density
    - Viscosity

# ML for Formulations

- Create ML models for mixtures and formulations containing multiple molecules

- Identify most import features and descriptors

- Design/optimize new formulations with novel compositions and chemistries



**Scheduled for 24-1**

*The technical features and projected timeline presented on this slide is for discussion purposes only. Such planned or potential capabilities are subject to change at any time.*

# Active Learning and Genetic Optimization

# Active Learning OptoElectronics Multi-Parameter Optimization (MPO)

# Active Learning Workflow for OptoElectronics

# Optoelectronic Genetic Optimization



AutoQSAR consensus prediction for experimental $\Delta E_{ST}$ dataset

# Machine Learning Forcefields

# Neural Network Potentials (NNPs)



Achieving QM accuracy at the cost of classical forcefields is an exciting prospect for neural network potentials to accelerate design of next-generation materials

# Our First NN Model: Schrödinger-ANI (SANI)



AEV table:

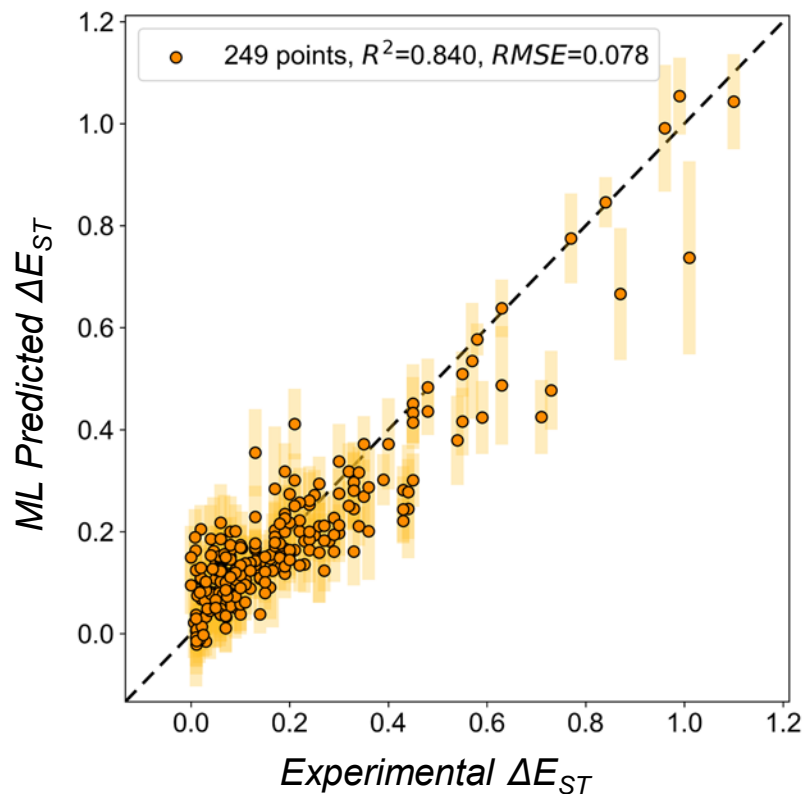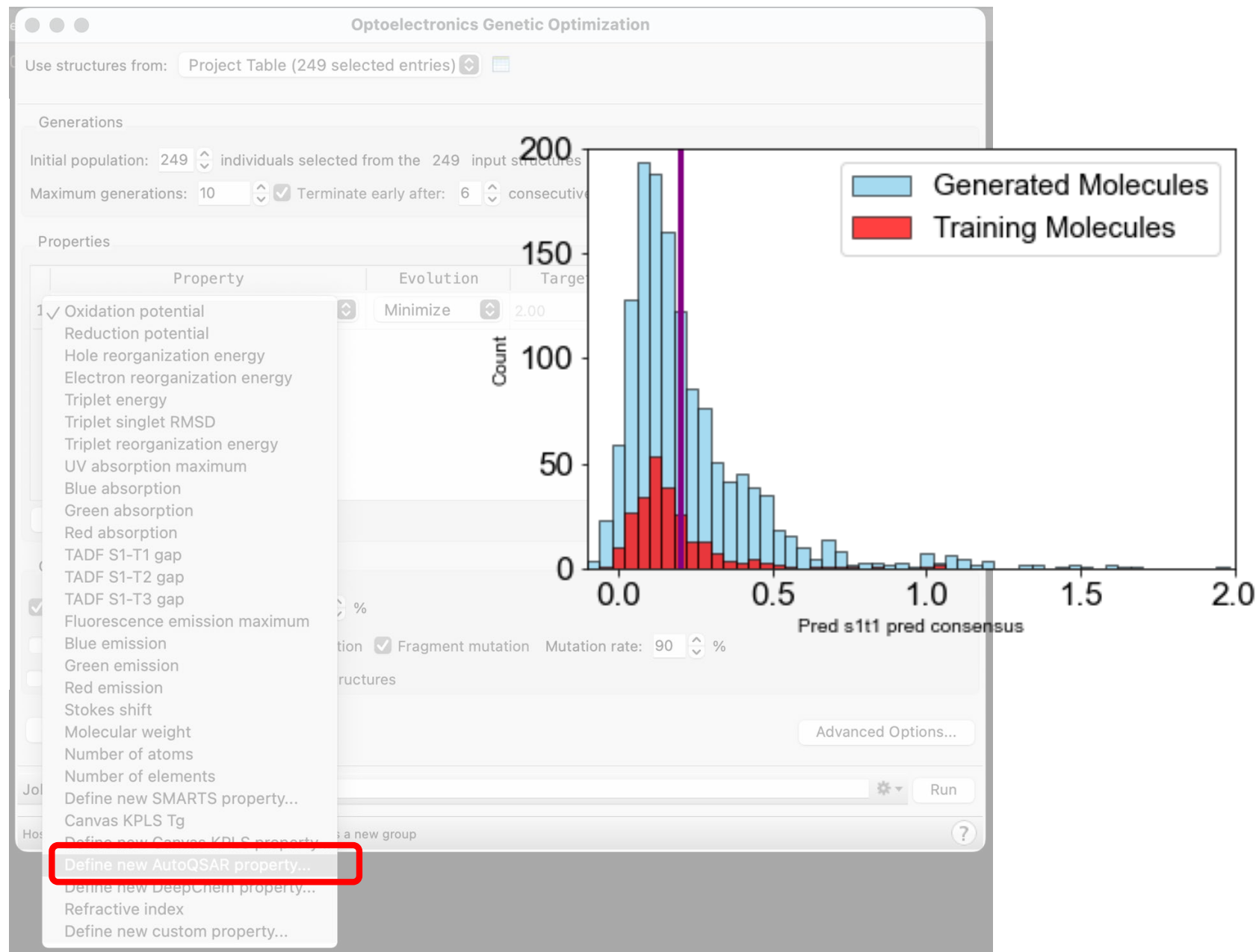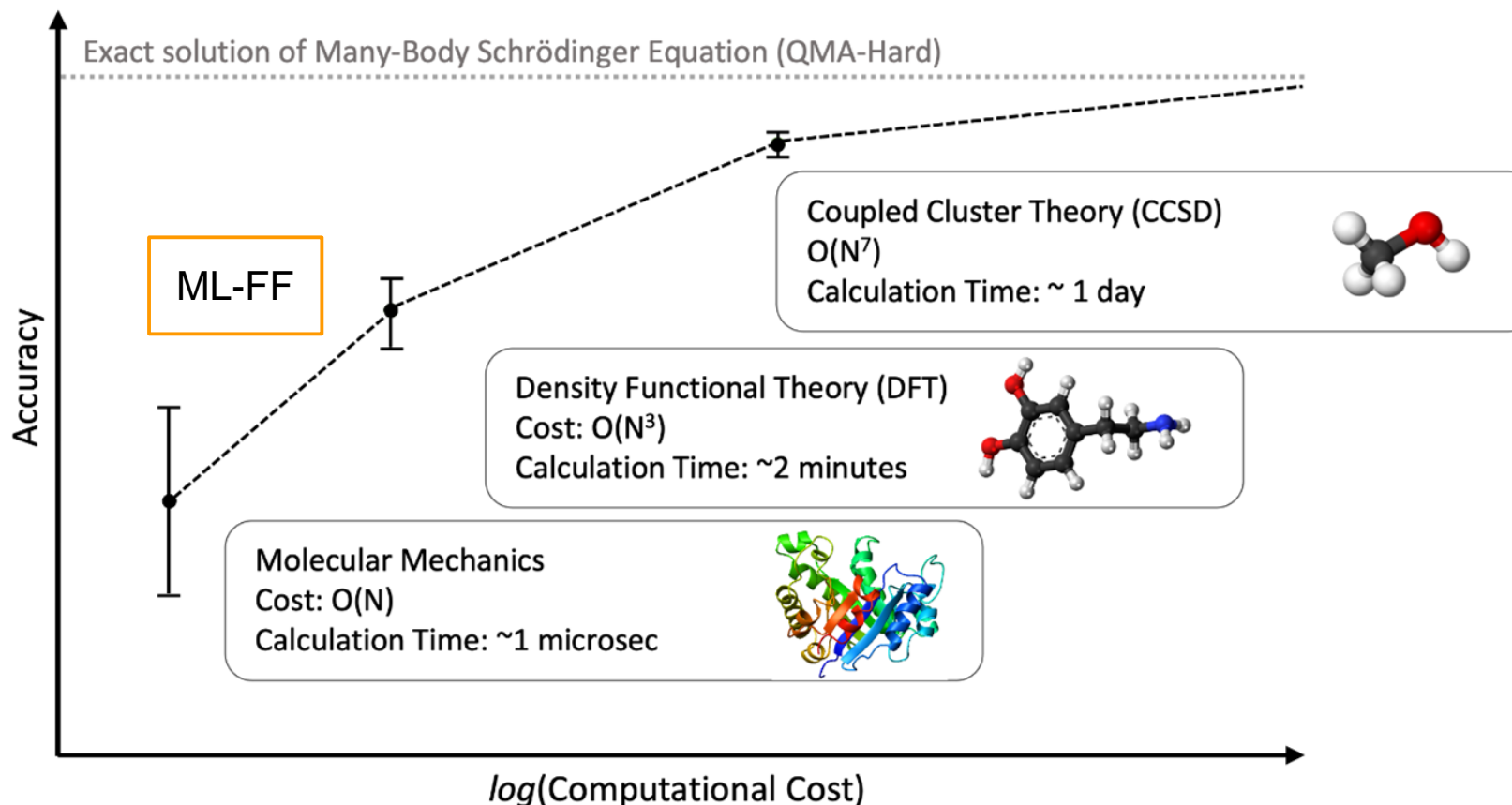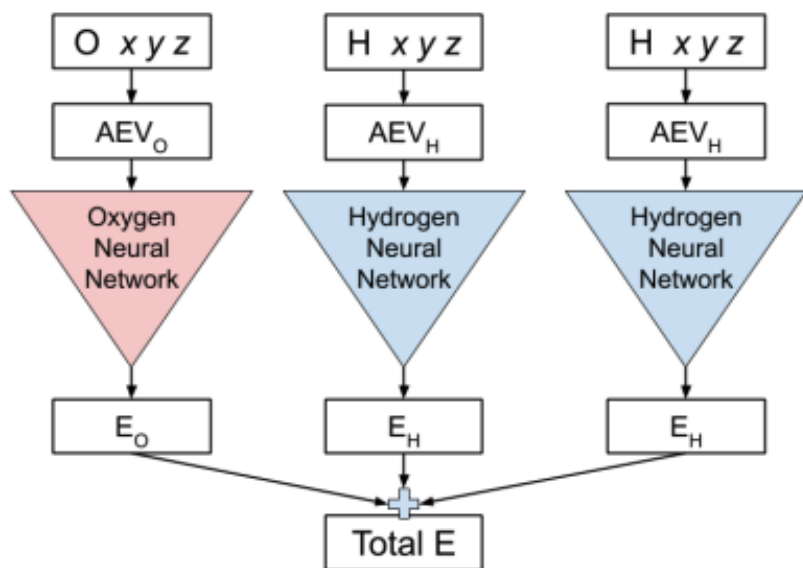| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| $AEV_H$ | $r_{HH}$ | $r_{HC}$ | $r_{HN}$ | $r_{HO}$ | $\theta_{HHH}$ | $\theta_{HHC}$ | $\theta_{HHN}$ | $\theta_{HHO}$ | $\theta_{CHH}$ | ... |
| $AEV_C$ | $r_{CH}$ | $r_{CC}$ | $r_{CN}$ | $r_{CO}$ | $\theta_{HCH}$ | $\theta_{HCC}$ | $\theta_{HCN}$ | $\theta_{HCO}$ | $\theta_{CCH}$ | ... |
| $AEV_N$ | $r_{NH}$ | $r_{NC}$ | $r_{NN}$ | $r_{NO}$ | $\theta_{HNH}$ | $\theta_{HNC}$ | $\theta_{HNN}$ | $\theta_{HNO}$ | $\theta_{CNH}$ | ... |
| $AEV_O$ | $r_{OH}$ | $r_{OC}$ | $r_{ON}$ | $r_{OO}$ | $\theta_{HOH}$ | $\theta_{HOC}$ | $\theta_{HON}$ | $\theta_{HOO}$ | $\theta_{COH}$ | ... |

*Stevenson et. al., arXiv, 1912.05079 (2019)*

❏ SANI is extension of ANI[1]-family of NN potential[2]

❏ Supports 8 elements covering 94% druglike molecules in ChEMBL[3]

❏ Inputs are cartesian coordinates and element type for each atom

❏ Each element has a separate NN learning mapping from features to energies

❏ Trained to DFT energies (wB97X/6-31G(d))

❏ Limitations: Neglect long-range effects, no information about charge state, charge distribution

[1]Smith et. al.,Chem. Sci., 8, 3192-3203 (2017)
[2]Behler et. al., Phys. Rev. Lett., 98, 146401 (2007)
[3]Gaulton et. al., Nucleic Acids Res., 45, D945-D954 (2016)

Schrödinger

# QRNN: Charge-Recursive Neural Network



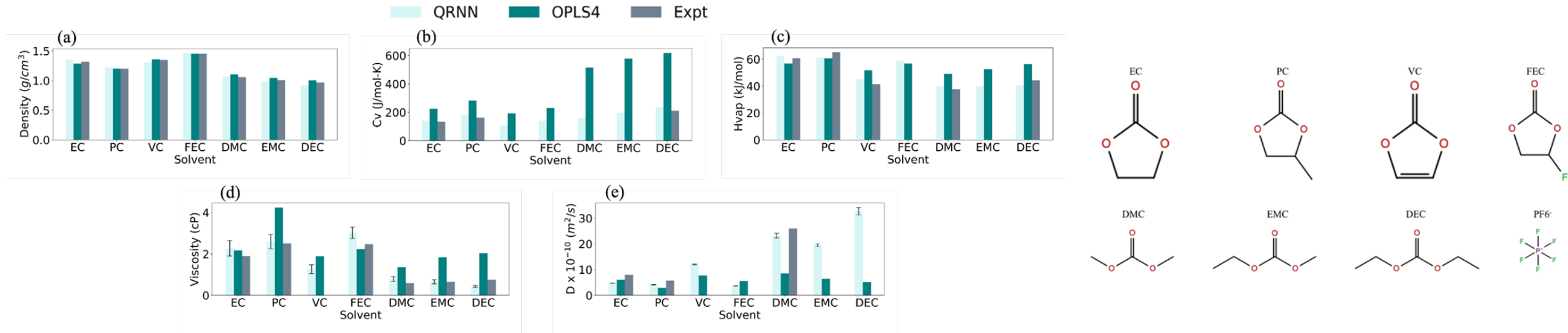Jacobson et. al., J. Chem. Theory Comput. 18, 2354-2366 (2022)

**Predicts energy, atomic charges, dipole and atomic forces**

❑ Extension of SANI to provide support for ionic systems

❑ Involves recursive charge correction

❑ Predicted atomic charges added as features to the NN

❑ Include charge dependent AEVs to learn radial charge distribution

❑ Empirical dispersion correction and coulomb interaction using the predicted charges added to energy

# Bulk Properties of Liquid Electrolytes



- ❏ ML-FF computed bulk properties are compared to experiments and OPLS4

- ❏ Excellent property predictions for all electrolyte components

- ❏ Significant improvement in prediction of diffusivity and viscosity compared to OPLS4

*The technical features and projected timeline presented on this slide is for discussion purposes only. Such planned or potential capabilities are subject to change at any time.*
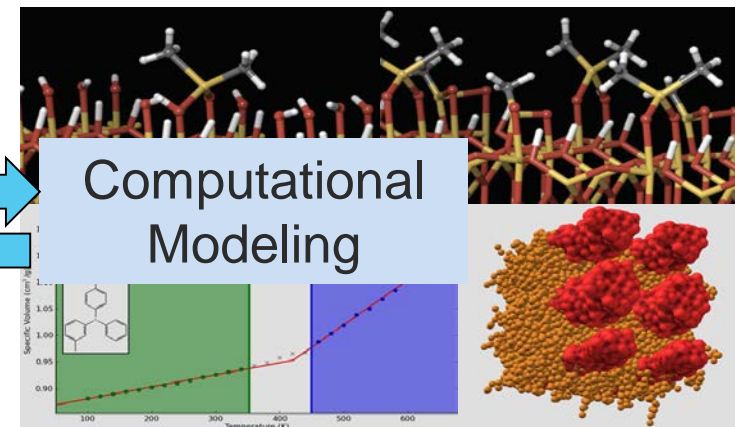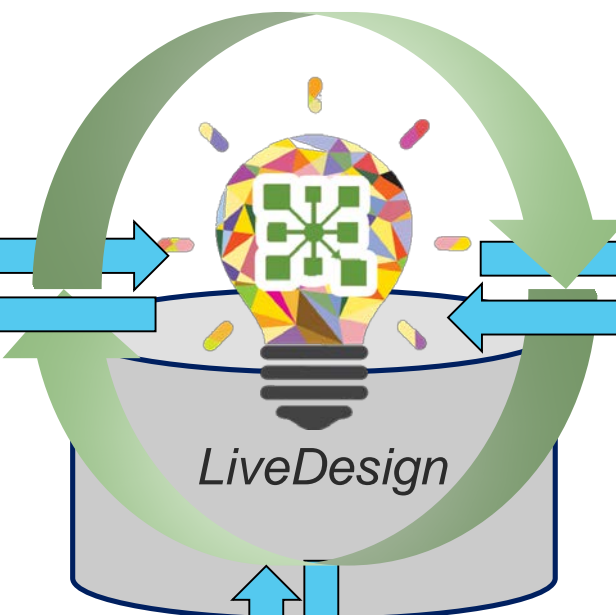
# Enterprise Informatics
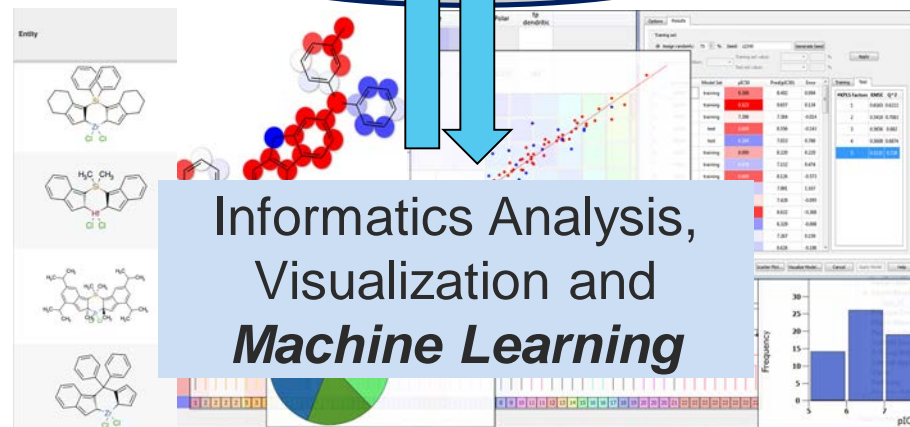
# Schrödinger's Informatics Platform - LiveDesign®



**Experimental data**

**Computational Modeling**

*LiveDesign*

**Informatics Analysis, Visualization and *Machine Learning***

- Web-based: Instantly LIVE to all users
- Scalable: Performant for global sized organizations
- Informatics: Visualization and Analysis
- Central Platform

- Easy agnostic access to expert computational tools
  – Machine learning
  – Advanced QM properties
- Execute modeling jobs, analyze results alongside all other data

*Schrödinger's core values of modeling supported DESIGN*

**Schrödinger**

# Suitable for Diverse Materials and Data Types



Polymers and co-polymers

Molecules

Experimental data

Formulations

Organometallics

# Summary

- Customized featurization based on chemical domain knowledge is critical in developing machine learning models

- High throughput physics-based modeling (QM, MD) provides various advantages in enhancing machine learning technology

- Our machine learning technology has been successfully applied to a wide range of materials systems and can be easily adapted to experiment design (i.e experimental design)

- Web-based materials informatics platform (LiveDesign) enables data digitization with advanced data analysis/visualization and machine learning technology

- ML Forcefields offer MD simulations with DFT-level accuracy

- Schrödinger's technology is to **empower** users and increase efficiency and productivity

Schrödinger

**Schrödinger**

**Thank you**

Feel free to reach out to me at:
chandras@schrodinger.com