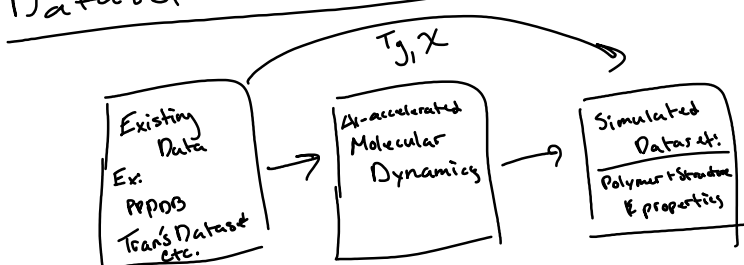


Pipeline Ideation

Sunday, June 26, 2022 2:01 PM

Dataset Creation:



ML Models for Property Prediction



ML Models to Consider

- Gaussian Process Regression
 - + works well on smaller datasets (less MD simulation)
 - + utilized in Polymer Genome (existing ML-based PPP platform)
 - + uncertainty measurements
 - not sparse (use all features), less efficiency when features exceed few dozen
- Deep Reinforcement Learning
 - + Molecular Generation model
 - Not as transferable to property prediction (DNN generally don't perform as well for PPP?)
- Random Forest
 - + Concatenating features and select for most important features after training (i.e. MF + QSPR descriptors + ME)
- CNN
 - + identified as possible competitive model (CNN-Representative ME)
 - 1D CNN exists
- Molecular Generation by Inverse Design (MGenID; note by IBM)
 - Not necessarily ML, but whatever
 - Feature Search by Particle Swarm Organization for most influential substructures from existing molecules.
 - Mainly algorithmic.
- Retrosynthesis by Molecular Transformer (IBM)
 - Utilized in 10M RXN for Chemistry
 - requires significant chemical reaction dataset
 - SMILES $\xrightarrow{2,3,4}$ (1 guess)
- GeoMol
 - + Provides valid conformers
 - + ties in with RDKit
 - + affects descriptors & vectorization

Guiding Objectives:

Unity / Low-Budget:

Tool for gaining insight through Machine

Autonomous / High-Budget:

Autonomous tool for while optimization

Pipeline Necessities:

- Polymer Property Prediction
- Polymer Generation op
- Way to screen generated

Molecular Dynamics:

- Utilize TorchMD for
- ↳ Uses PyTorch (i)
- ↳ Designed for prot

ML Backends:

- Scikit-Learn for uti
- PyTorch for TorchMD

* utilize to building b1 for DeepRL

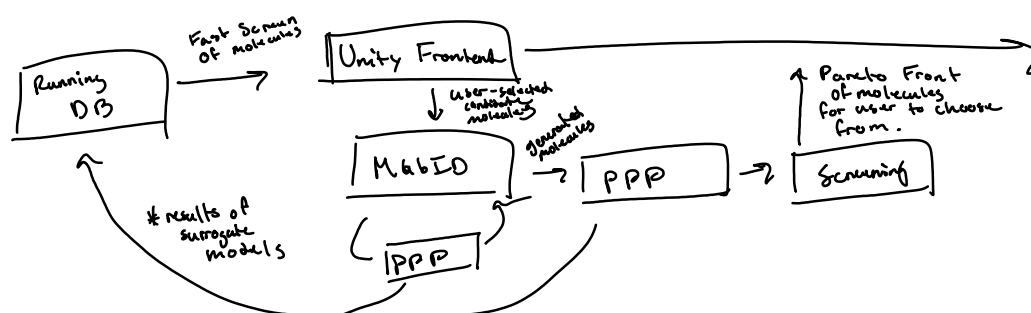
Efficient multi-objective molecular optimization in continuous latent space†

Robin Winter ^{1,2,3}, Floriane Montanari ^{1,2,3}, Andreas Steffen ^{1,2,3}, Hans Briem ^{1,2,3}, Frank Noé ^{1,2,3} and Dj

UNITY PIPELINE

Training Dataset:

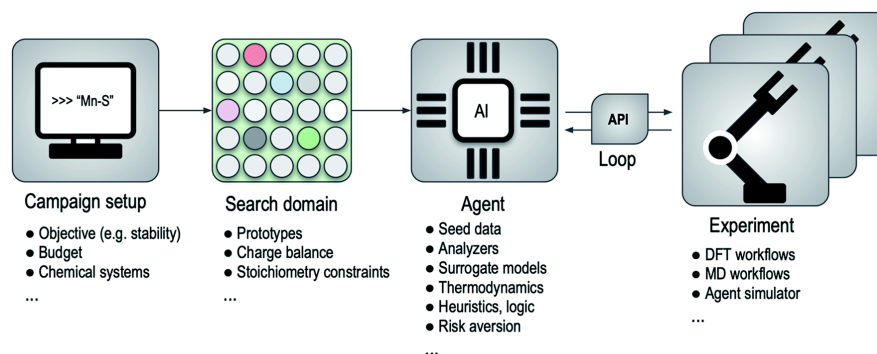
- DO NOT generate new molecules for training
- Minimum available dataset, i.e. implement w/ least amount of simulations



AUTOMATED PIPELINE

1) Computational Autonomy for Materials Discovery

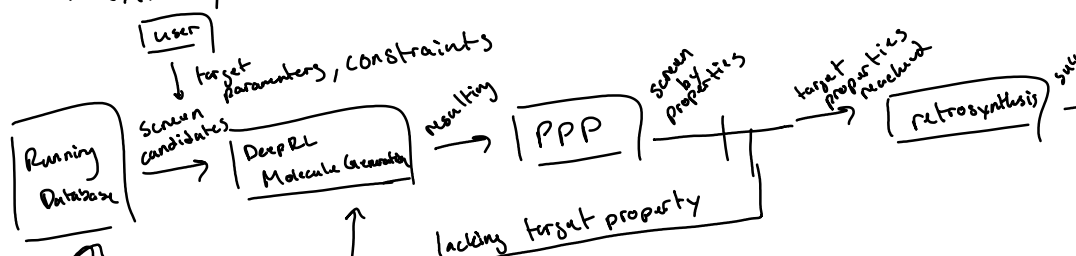
- see Montoya et al
- code on github; open-source framework



2) My "Brash" Pipeline

Training Dataset

- "Randomly" generate new molecules for training to fill in the gaps
- Utilize Geo Mol to enhance featurization.
- Maximally available dataset



add to DB and continue until criterions
are met
or molecule(s) is Pareto optimal