# Machine Learning to Predict Experimental Protein-Ligand Complexes THE GEORGE

Hyunji Kim<sup>1</sup>, Sarah Walworth<sup>2</sup>, Kenny Merz<sup>3</sup>, Jun Pei<sup>3</sup>, Lin Song<sup>3</sup>, Zheng Zheng<sup>3</sup>

# OBJECTIVE

Traditional scoring methods to determine correct poses for protein-ligand binding are generally around 60% accurate. Our goal was to use random forest machine learning to optimize the ability to predict ligand poses that are close to the native crystal structure of the protein-ligand complex. One major application of our method is **drug design**. It will allow "designers" to find molecules that could dock similarly to the native crystal structure.

## METHOD

#### **Data Generation:**

Given 766 protein-ligand complexes, we generated ligand decoys (up to 100 per protein) using Schrodinger Glide software.

#### **GARF Potential Function:**

We considered approximate effects on energy using the GARF pairwise interatomic potential function.

#### 3. <u>Random Forest (RF)</u>

We generated our RF model using the Scikit-Lean tool. **5-fold cross validation** was applied to prevent overfitting. For every simulation, we randomly selected 70% of the data as the training set and 30% as the test set.

#### *Classification:*

 $Rank_{native} < Rank_{ligand_{decoy}} \Longrightarrow Class 0$  $Rank_{native} > Rank_{ligand_{decov}} \Longrightarrow Class 1$ 

#### 4. <u>Scoring</u>

We used the Cambridge Crystallographic Data Center (CCDC) GOLD protein-ligand docking software to generate the Astex Statistical **Potential(ASP)** and **Chemscore** scoring functions to validate our RF model.

#### 5. Post-Processing

We performed **17 grid searches** to narrow down the RF parameters. From the grid search, we **selected the best 6 parameters** and ran 12 independent simulations for each of the 6 parameter combinations to identify the best parameter for our RF model.



Figure 3: Graphical representation of the Random forest model <sup>1</sup>

# <sup>1</sup>George Washington University,<sup>2</sup>University of Colorado Boulder, <sup>3</sup>Michigan State University

Table 1: Final RF Parameter (Choose Parameter 6 from figure 2)

<sup>1</sup>T. (2018, May 01). Machine Learning for Beginners, Part 10: Random Forest. Retrieved from https://thedatalass.com/2018/04/17/random-forest/







VALIDATION **Scoring Function Comparison** 1.00 0.80 0.40 0.20 0.00 **RF** Test Chemscore ASP 0.657 0.6 0.923

Figure 3: Comparison of RF results against the results of two traditional scoring functions (Chemscore and ASP)

#### Chemscore:

**Empirical scoring function**: Regression based with coefficients based on experimental data, which accounted for physical factors that affect docking.

#### **Astex Statistical Potential (ASP):**

- Atom to atom potential function using the Worldwide Protein Data Bank.
- Considered frequency and potentials: Expected number of interactions of atoms in a defined radius.

#### CONCLUSION

#### **Overall Results:**

• Our results have shown that our random forest machine learning model is **significantly more accurate** in predicting ligand poses similar to the native crystal structure of a protein-ligand complex than two traditional scoring functions.

#### **Future Work:**

- For further validation, we plan to **test our model with a larger data set** using the sets of decoys that have been generated in **the** Database of Useful Decoys: Enhanced (DUD-E).
- We also plan to use the GARF scoring function as an accuracy comparison.

### REFERENCES







Research Experience for Student

We acknowledge support from the MSU ACRES REU program, which is supported by the National Science Foundation through grant ACI-1560168.

