

Charlie Li: Data Prep, Transformers, Feature Importance

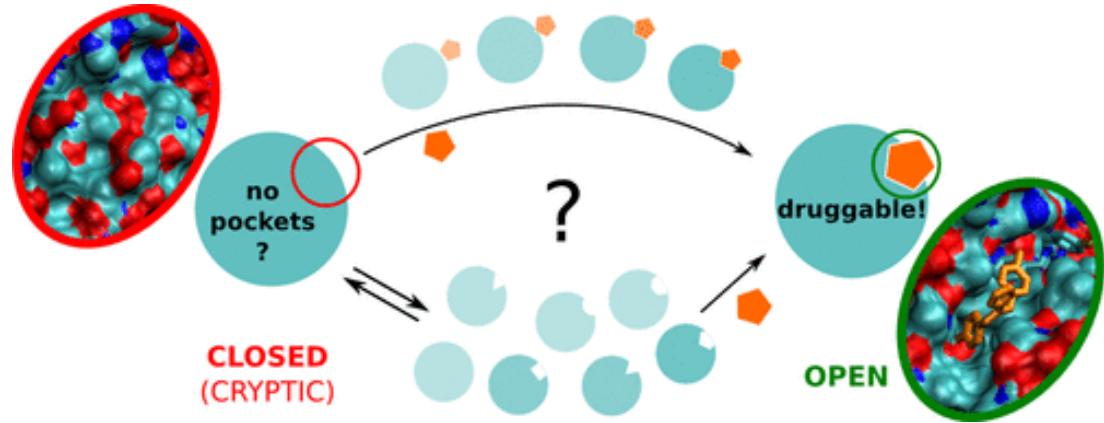
Matthew Rao: Data Prep, Single Models, Hybrid Models

Contents

- 1. Introduction and background information**
2. Dataset
3. Model Training
 - a. Transformer Encoder
 - b. Fusion Encoder
 - c. Single Models
 - d. Hybrid Models
4. Feature Importance
5. Outlook

Cryptic Pockets

- Solvent exposure only upon ligand binding²
- Potential drug targets



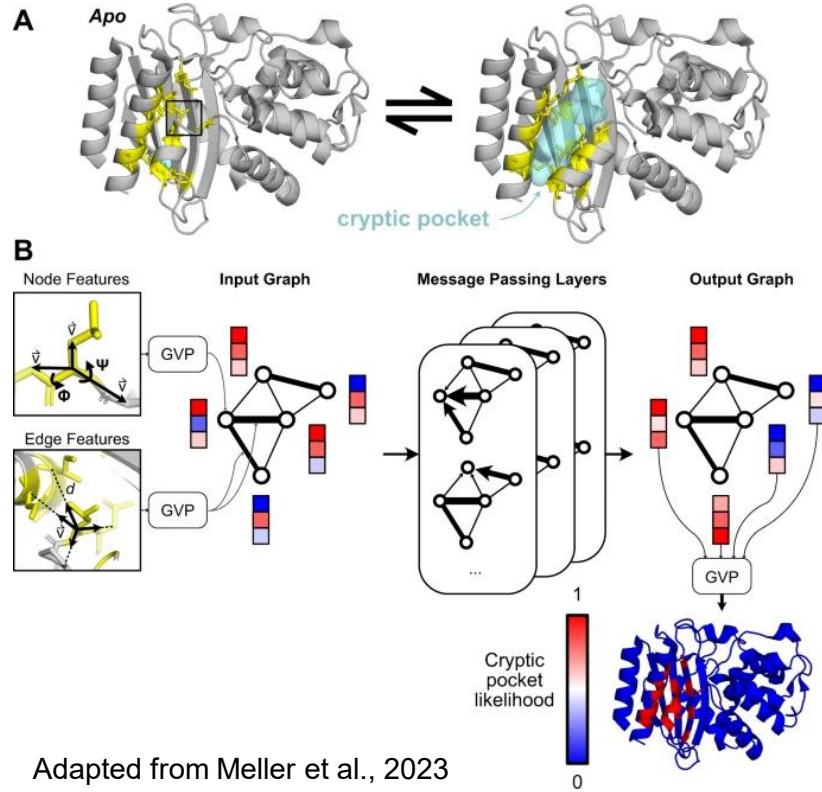
Challenges

- Dynamic property but static imaging
- Ground-state inaccessible: spontaneous closing in molecular dynamics simulation

Adapted from Oleinikovas et al., 2016

Past Work

- **CryptoSite**³
 - Simulation input on-the-fly
 - Trained on known examples
 - Run time: 1 day/protein
 - ROC-AUC: 0.83
- **PocketMiner**⁴
 - Graph neural network
 - Trained on simulations
 - Fast run time
 - ROC-AUC: 0.79 and 0.83



Adapted from Meller et al., 2023

Issues

- Representativeness of the starting point static structure
- Left out other relevant information
 - Evolutionary conservation, post translational modification, etc

Our Project

- Residue-level Crypticity Classification
- Sequence and Multiple Feature Categories
- Hybrid Model, ROC-AUC = 0.87

Contents

1. Introduction and background information

2. Dataset

3. Model Training

a. Transformer Encoder

b. Fusion Encoder

c. Single Models

d. Hybrid Models

4. Feature Importance

5. Outlook

First, Some Quick Definitions...

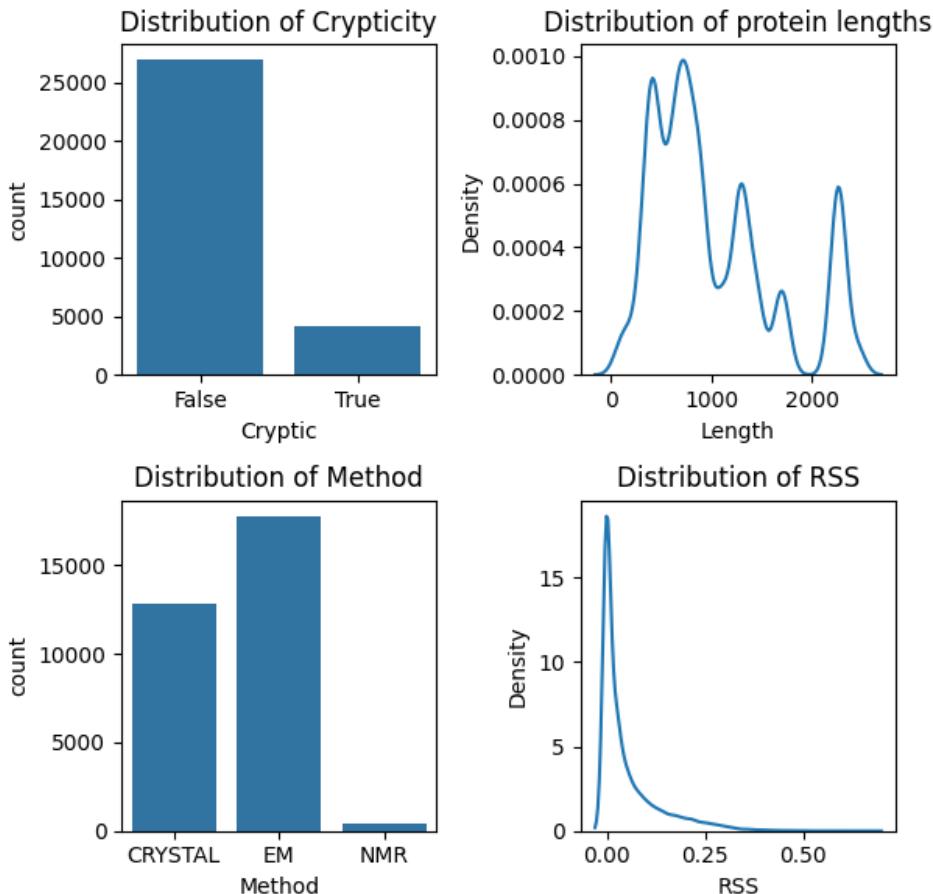
- ASA = Accessible Surface Area
- RSA = Relative Solvent Accessibility = $\frac{ASA_{observed}}{\max(ASA)} \times 100\%$
- Crypticity Criterion:

max(RSA)>25 && min(RSA)<25 && max(RSA) - min(RSA) > 20

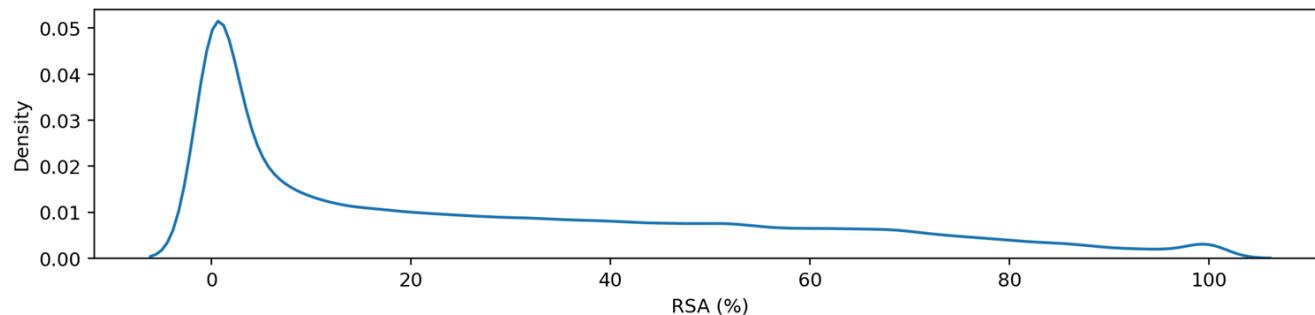
- Why 25 and 20? Arbitrary designation.

Data Sources

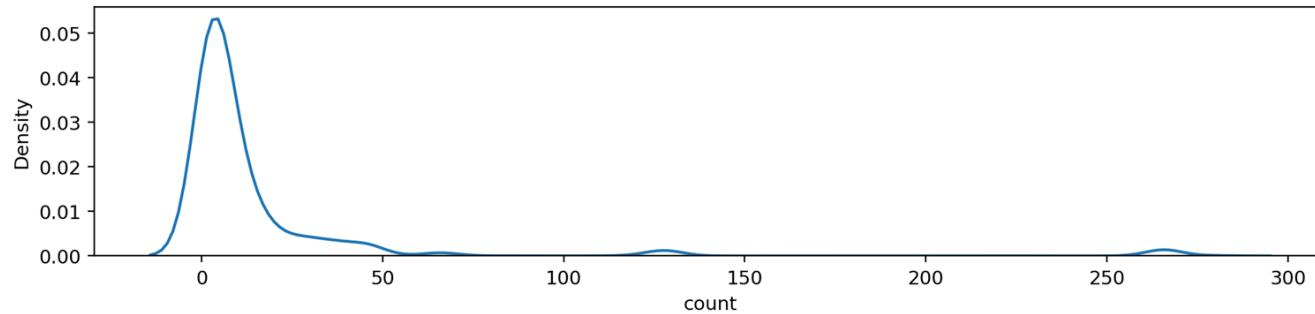
- **Protein Data Bank (PDB)**⁵: Sequence, Method, Annotations, residue RSA (calculated through a developed pipeline)
- **Uniprot**⁶: Protein names
- Amino Acid Chemical Properties
- **Aminode**⁷: Relative Substitution Score (RSS)
- **PhosphoSite**⁸: Post Translational Modifications (No Match)
- CSV Table
- 500,000 Observations
- After aggregation: 31,080 rows
- After dropping RSS mismatches: 27,844 rows



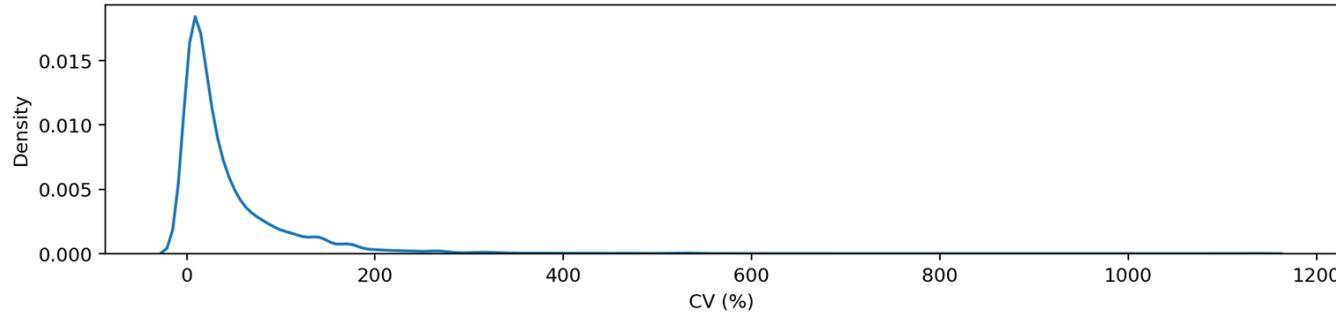
Distribution of RSA Across All Sites



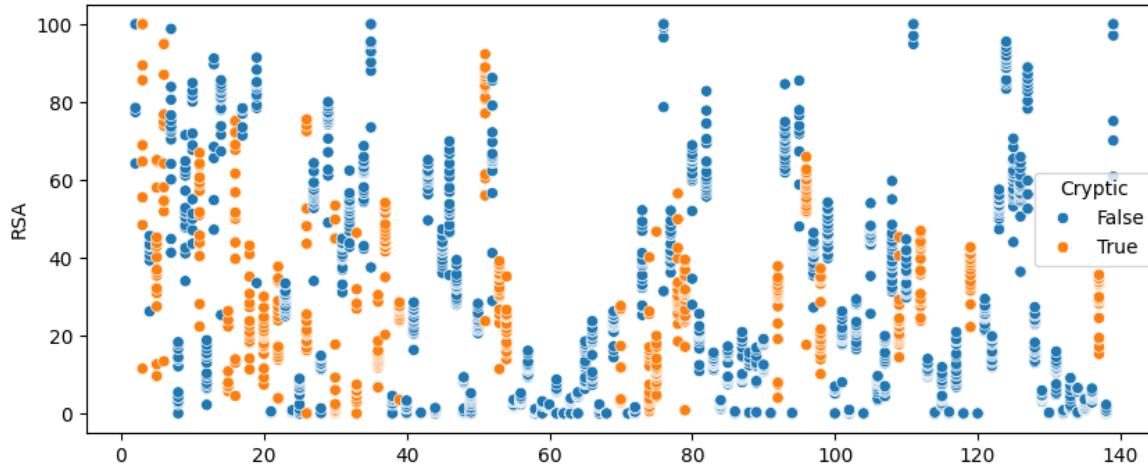
Distribution of Number of RSA Observations Per Site



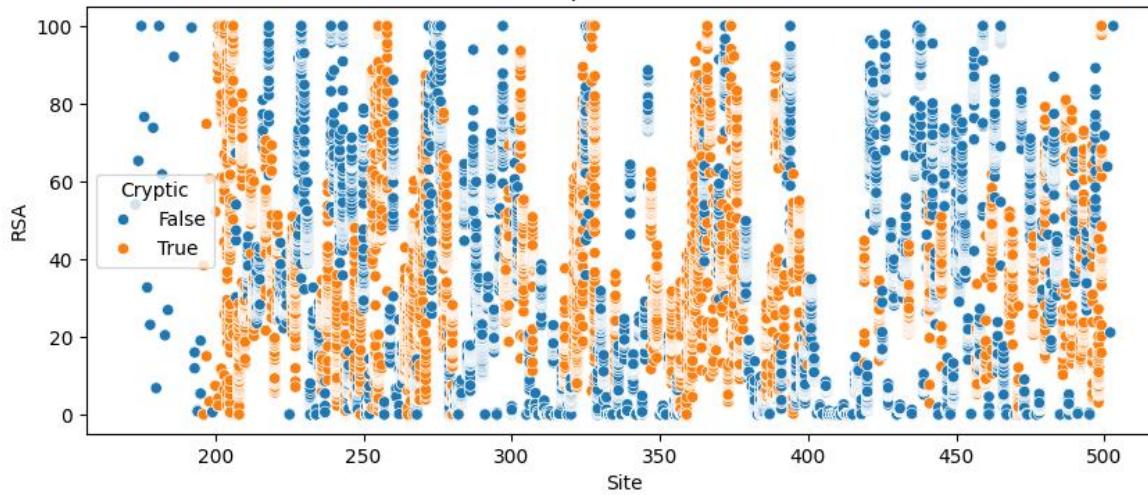
Distribution of RSA Coefficient of Variation Per Site



RSA Map of Gene ACOT13



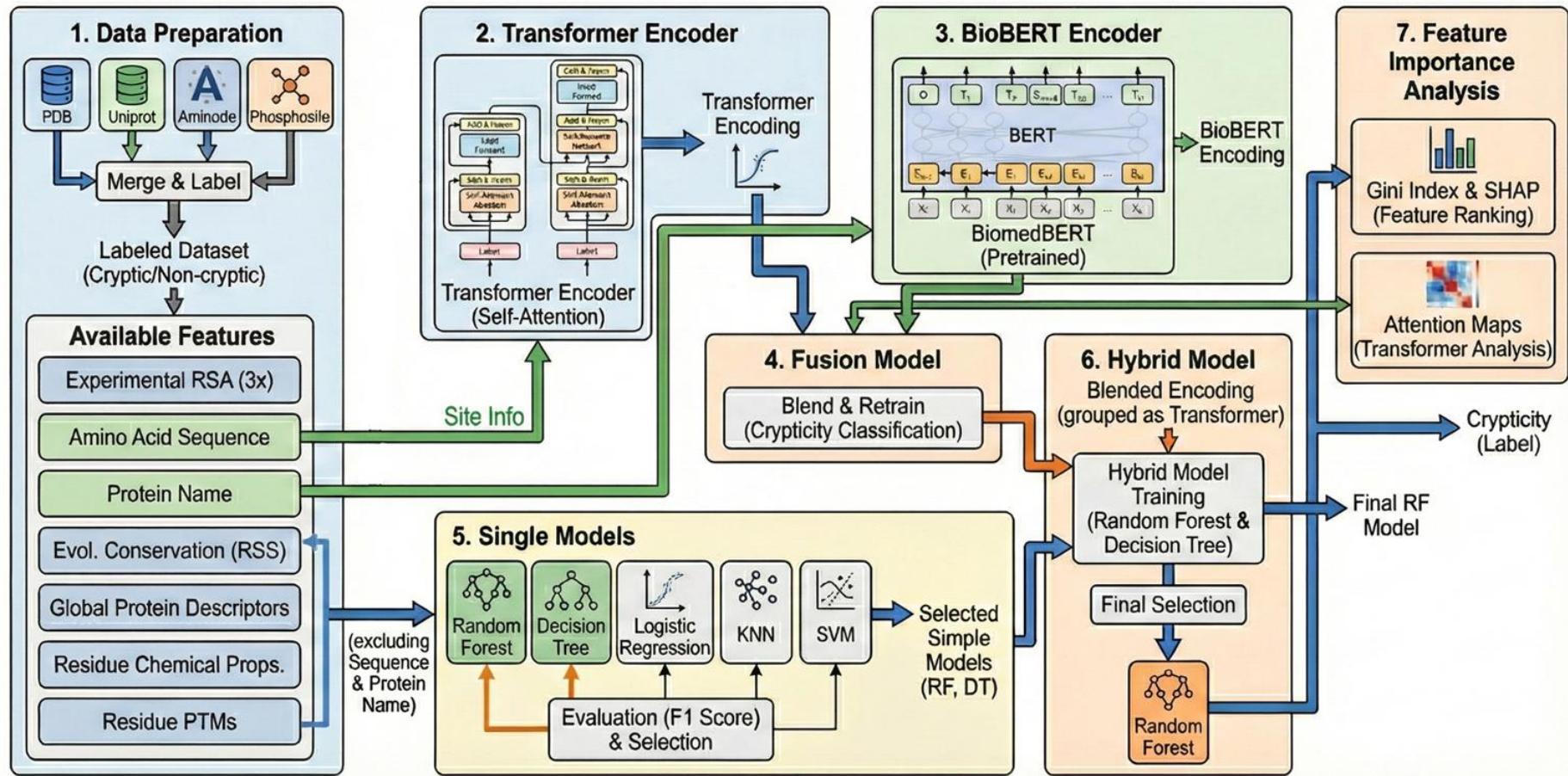
RSA Map of Gene ACVR1



Contents

1. Introduction and background information
2. Dataset
- 3. Model Training**
 - a. Transformer Encoder
 - b. Fusion Encoder
 - c. Single Models
 - d. Hybrid Models
4. Feature Importance
5. Outlook

Protein Cryptic Site Identification ML Workflow



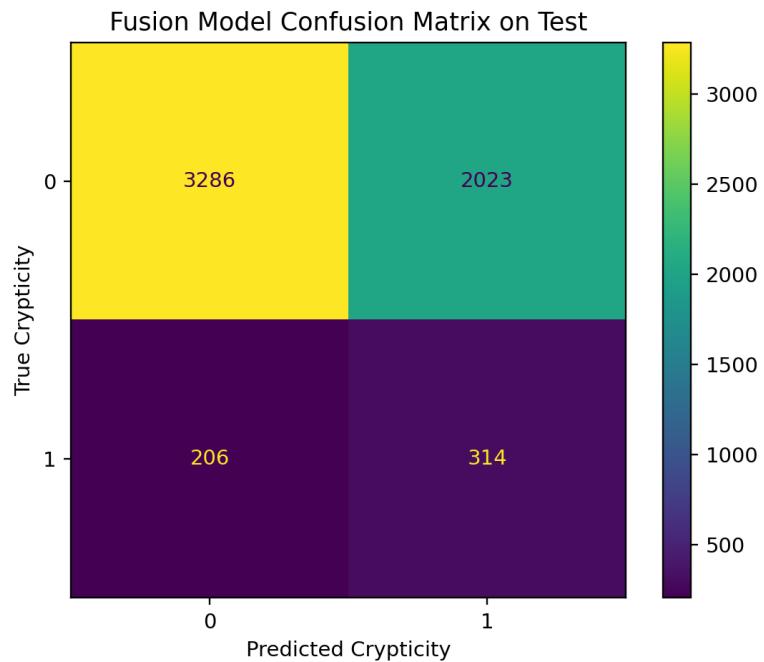
a. Transformer Encoder

```
| def forward(self, x, positions):
|     h = self.emb(x) + self.pe.unsqueeze(0)
|     key_padding = (x == self.padding_idx)
|     h2, attn = self.mha(h, h, h, key_padding_mask=key_padding, need_weights=True)
|     h = self.ln1(h + h2)
|     h2 = self.ff(h)
|     h = self.ln2(h + h2)
|     # extract encoding at target positions
|     batch_idx = torch.arange(x.size(0), device=x.device, dtype=int)
|     target_embedding = h[batch_idx, positions]
|     encoded = self.cls(target_embedding)
|     return encoded, target_embedding, attn
```

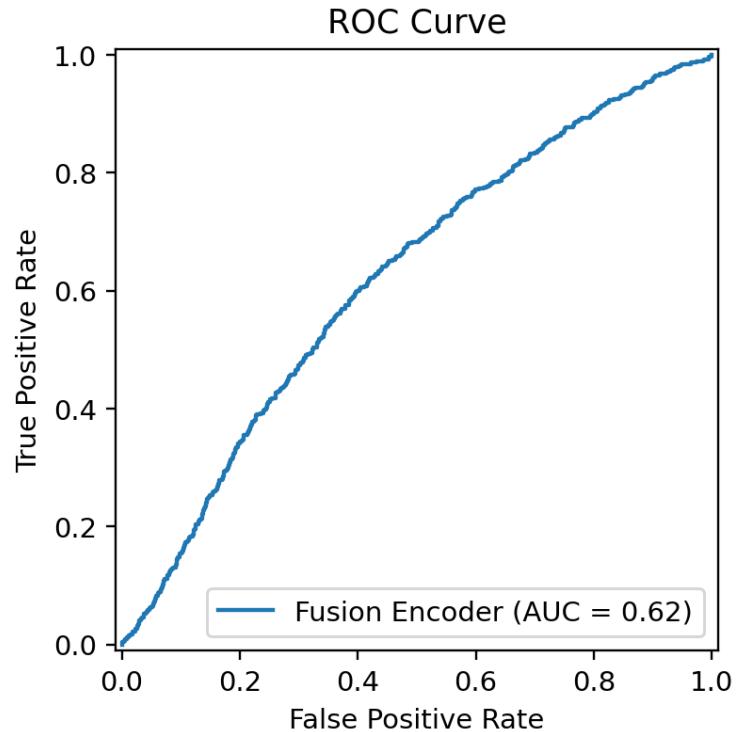
b. Fusion Model With Pre-trained BiomedBERT⁹

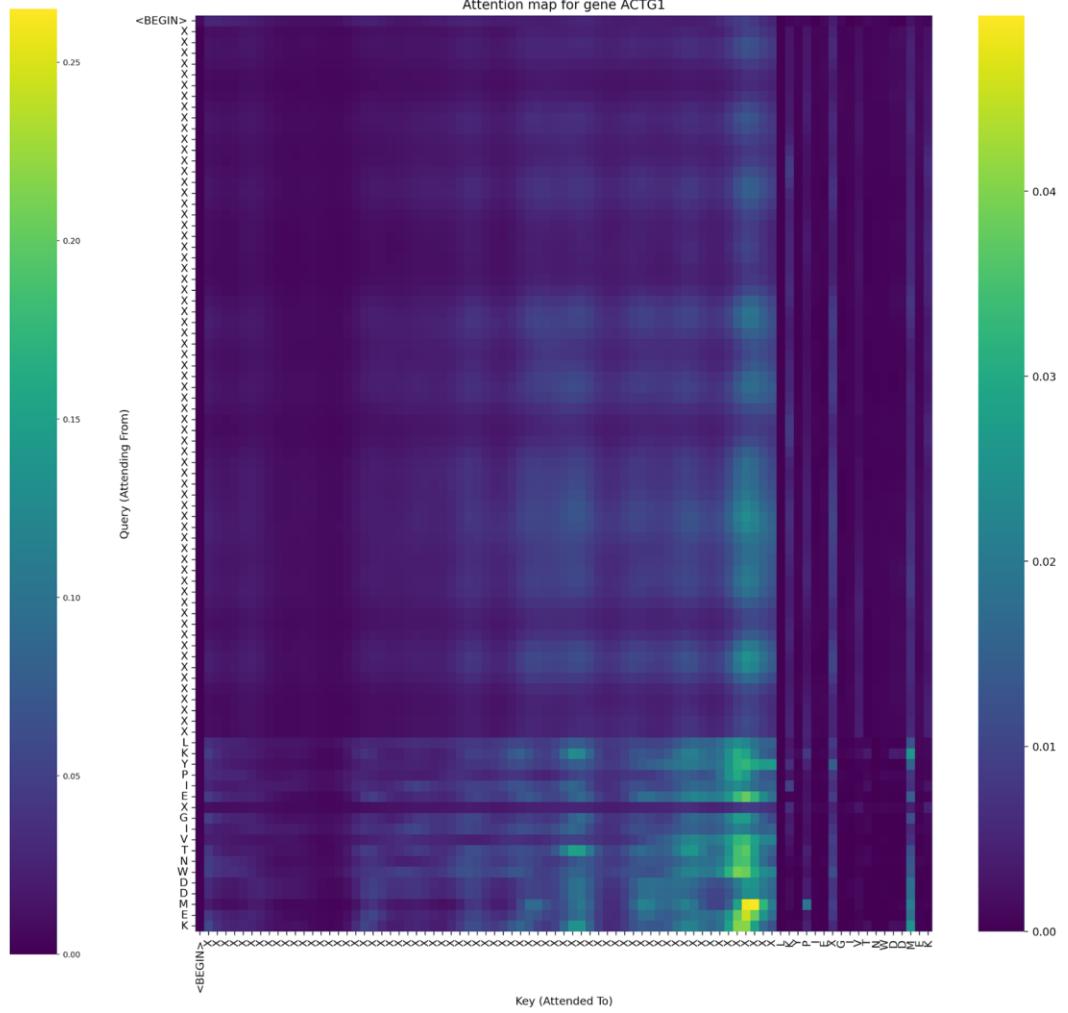
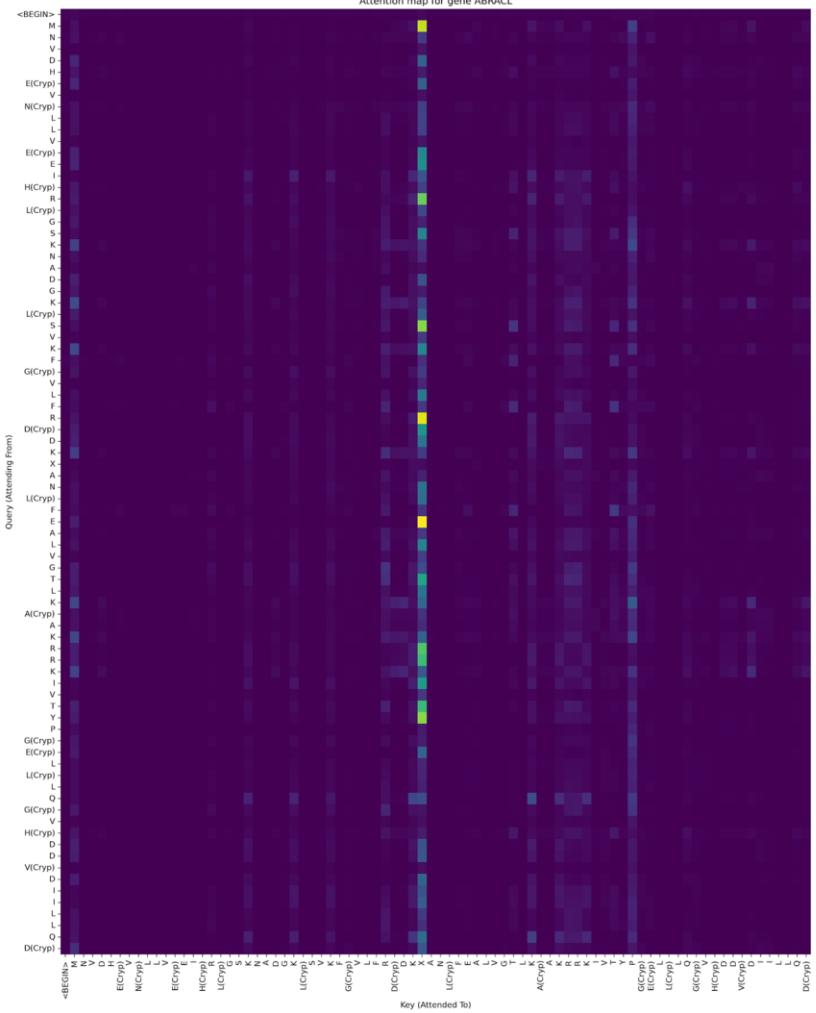
```
def forward(self, protein_encoding, site_encoding):
    # Modality Dropout (Training Only)
    if self.training:
        if torch.rand(1).item() < 0.3:
            protein_encoding = torch.zeros_like(protein_encoding)
        elif torch.rand(1).item() < 0.1:
            site_encoding = torch.zeros_like(site_encoding)
    mat = self.bn(protein_encoding)
    mat = self.ff1(mat)
    mat = self.ln1(mat)
    mat = self.ff2(mat)
    mat = site_encoding + 0.1*mat
    mat = self.ln2(mat)
    logit = self.cls(mat)
    return logit, mat
```

- Accuracy: 0.62
- F1 Score: 0.22
- ROC AUC: 0.621



Decent accuracy but poor F1 score





c. Single Model Analysis

- Utilized 5 Single Models
 - Random Forest
 - Logistic Regression
 - Decision Tree
 - KNN Classifier
 - Support Vector Machine
- Rank by F1 score

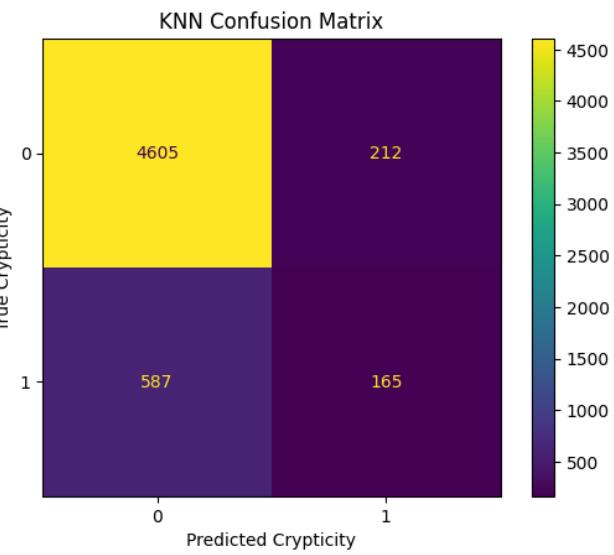
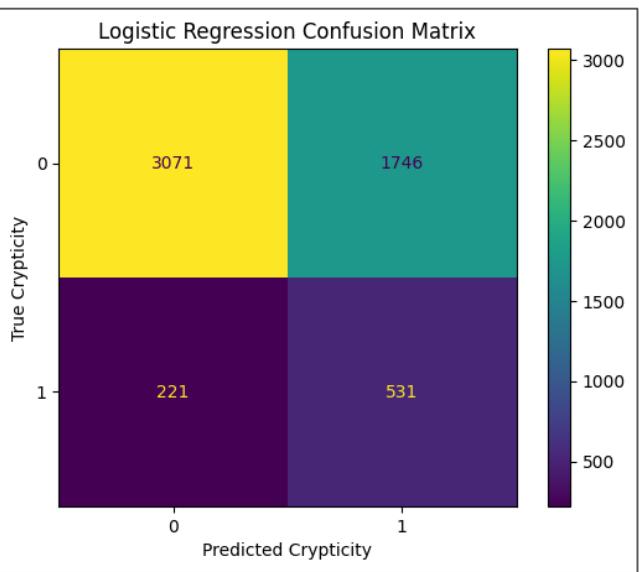
```
class_weight={0: 1, 1: 6.5} # it's 6.5 times more important to correctly classify a cryptic site
```

```
rf_model = RandomForestClassifier(  
    n_estimators=300,  
    min_samples_split=20,  
    min_samples_leaf=10,  
    class_weight=class_weight,  
    random_state=114514,  
    n_jobs=-1,  
    verbose=True  
)
```

```
▶ rf_model.fit(X_train, y_train)  
y_pred = rf_model.predict(X_test)  
  
print("Accuracy:", accuracy_score(y_test, y_pred))  
print("F1 Score:", f1_score(y_test, y_pred))  
  
... [Parallel(n_jobs=-1)]: Using backend ThreadingBackend with 2 concurrent workers.  
[Parallel(n_jobs=-1)]: Done 46 tasks      | elapsed:   1.7s  
[Parallel(n_jobs=-1)]: Done 196 tasks      | elapsed:  11.3s  
[Parallel(n_jobs=-1)]: Done 300 out of 300 | elapsed:  15.7s finished  
[Parallel(n_jobs=2)]: Using backend ThreadingBackend with 2 concurrent workers.  
[Parallel(n_jobs=2)]: Done 46 tasks      | elapsed:   0.1s  
Accuracy: 0.8125336685221763  
F1 Score: 0.4846989141164857  
[Parallel(n_jobs=2)]: Done 196 tasks      | elapsed:   0.2s  
[Parallel(n_jobs=2)]: Done 300 out of 300 | elapsed:   0.4s finished
```

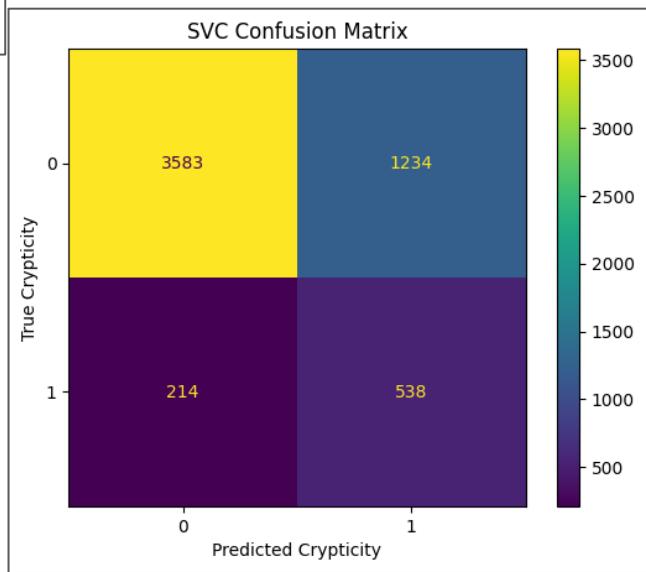
Logistic Regression

- Accuracy: 0.65
- F1 Score: 0.35



KNN

- Accuracy: 0.86
- F1 Score: 0.29

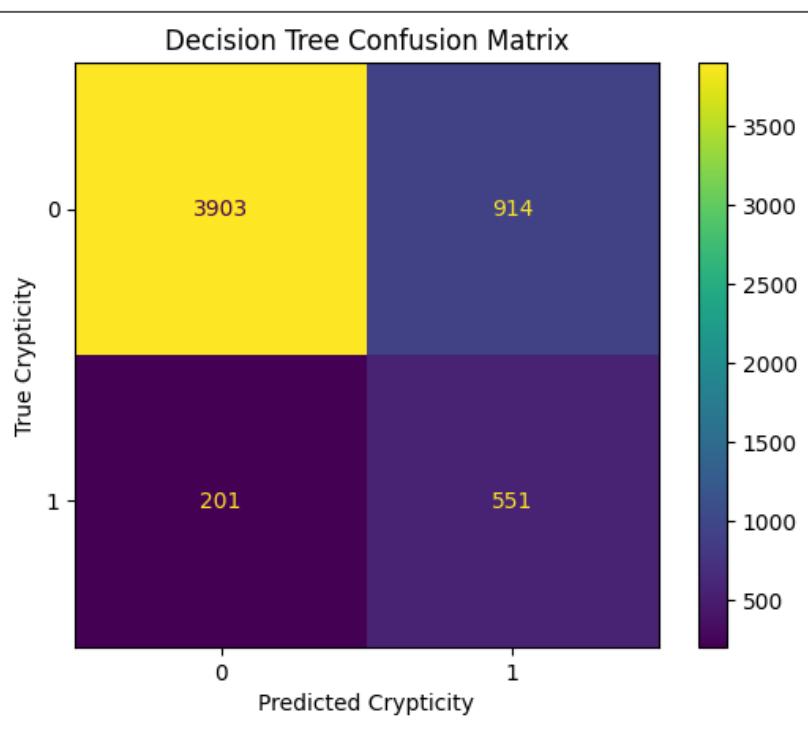


SVM

- Accuracy: 0.74
- F1 Score: 0.43

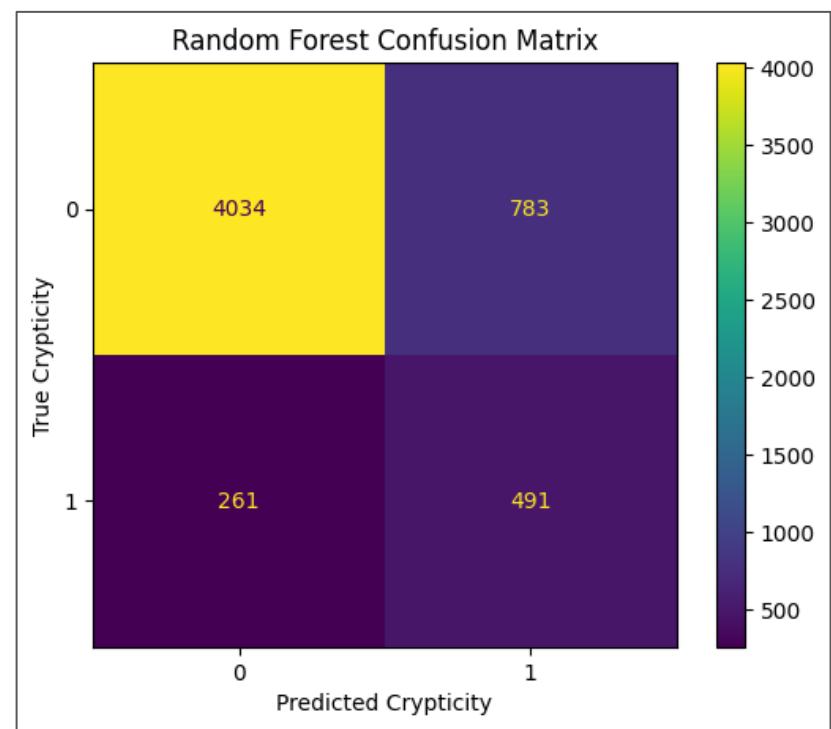
Decision Tree

- Accuracy: 0.80
- F1 Score: 0.50



Random Forest

- Accuracy: 0.81
- F1 Score: 0.48

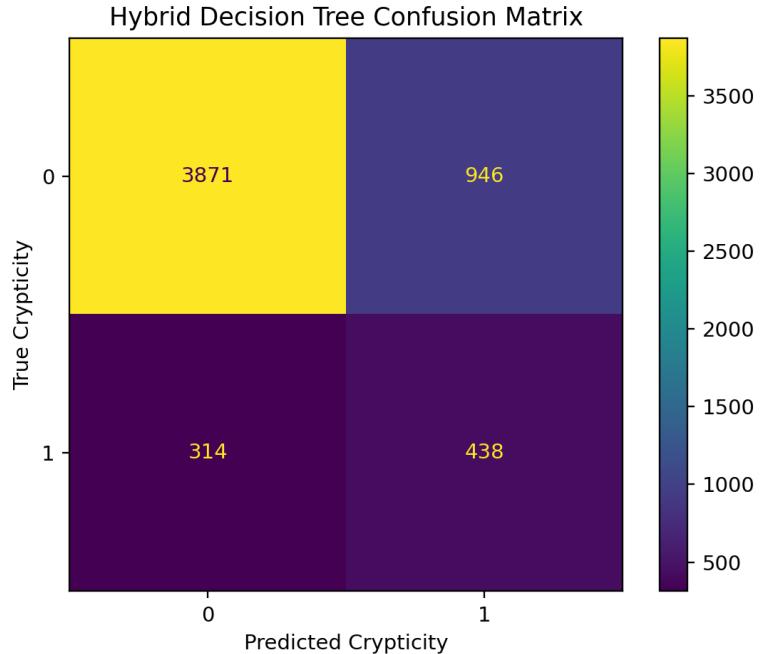


d. Hybrid Models

- Decision Tree, Random Forest
- Combined Features
 - Transformer encoding (64D vectors)
 - Three RSA random samples
 - Evolutionary conservation (RSS)
 - Global protein descriptors
 - Residue chemical properties
 - Residue post-translational modifications

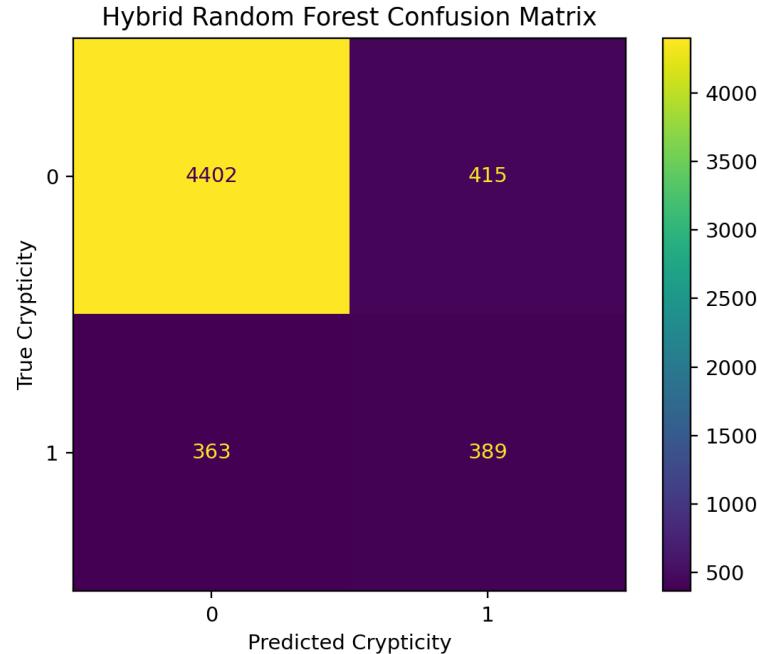
Hybrid Decision Tree

- **Accuracy: 0.77**
- **F1 Score: 0.41**
- **Performance dropped!**



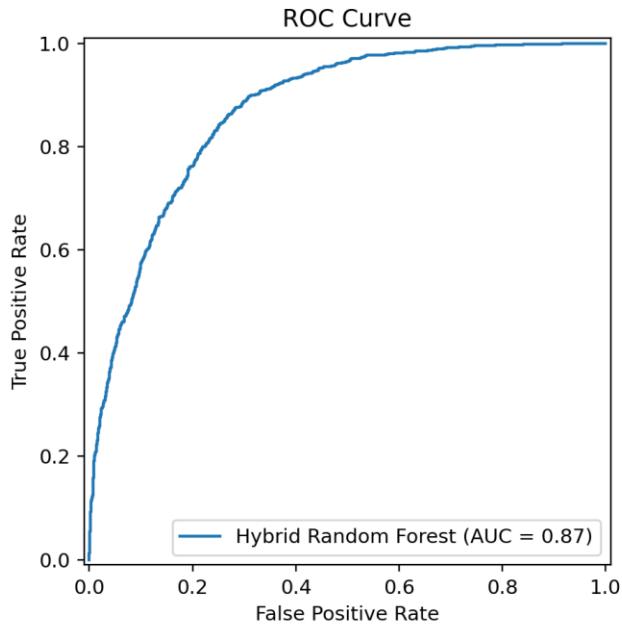
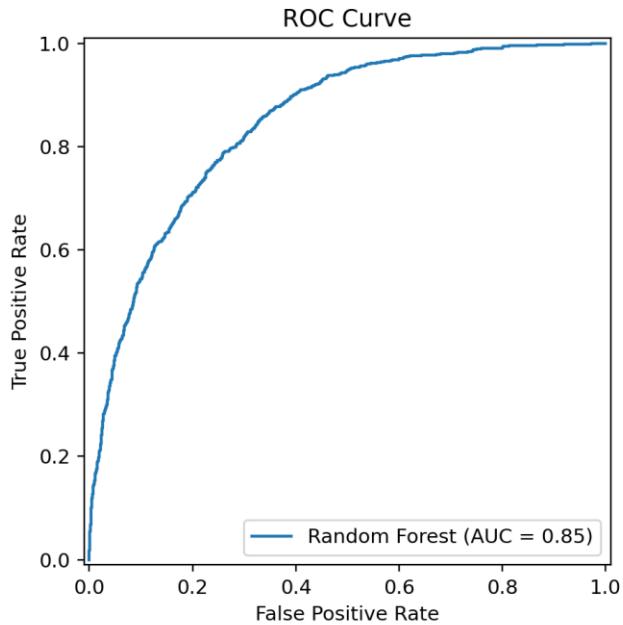
Hybrid Random Forest

- **Accuracy: 0.86**
- **F1 Score: 0.50**
- **More stable!**



ROC AUC before and after

Slight improvement!



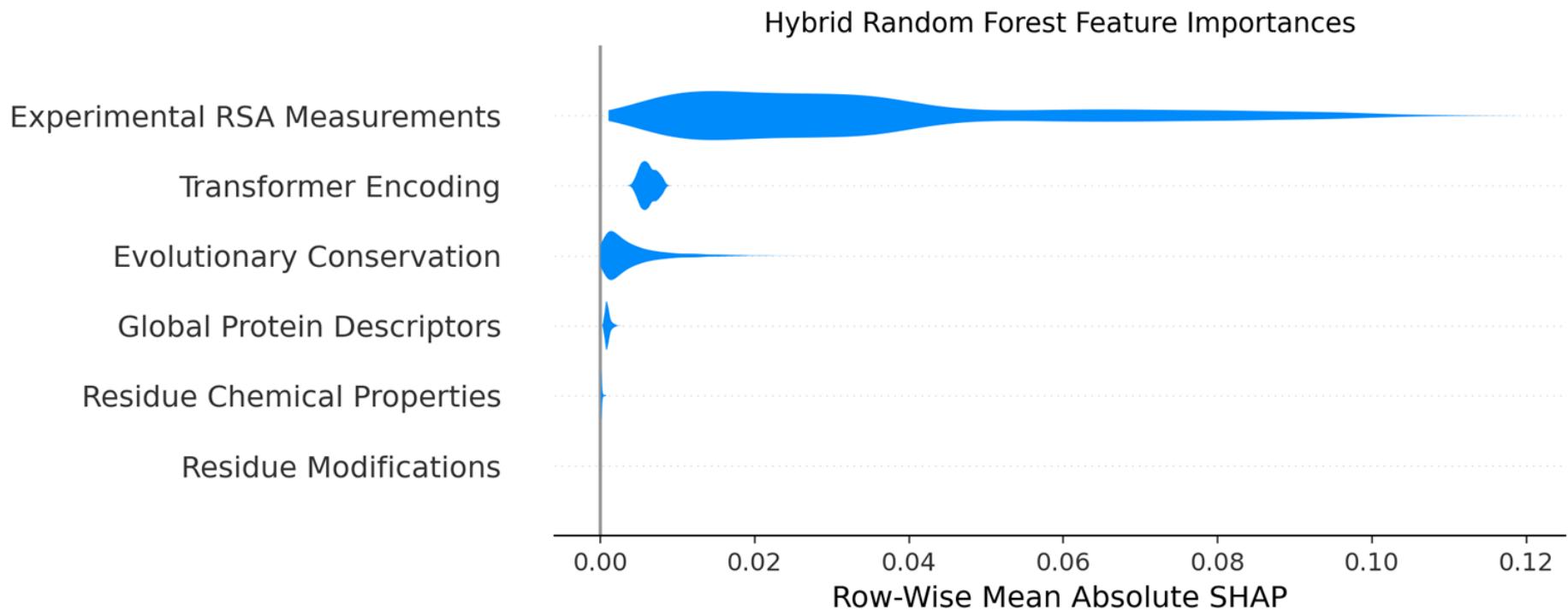
Contents

1. Introduction and background information
2. Dataset
3. Model Training
 - a. Transformer Encoder
 - b. Fusion Encoder
 - c. Single Models
 - d. Hybrid Models
- 4. Feature Importance**
5. Outlook

SHAP: SHapley Additive exPlanations

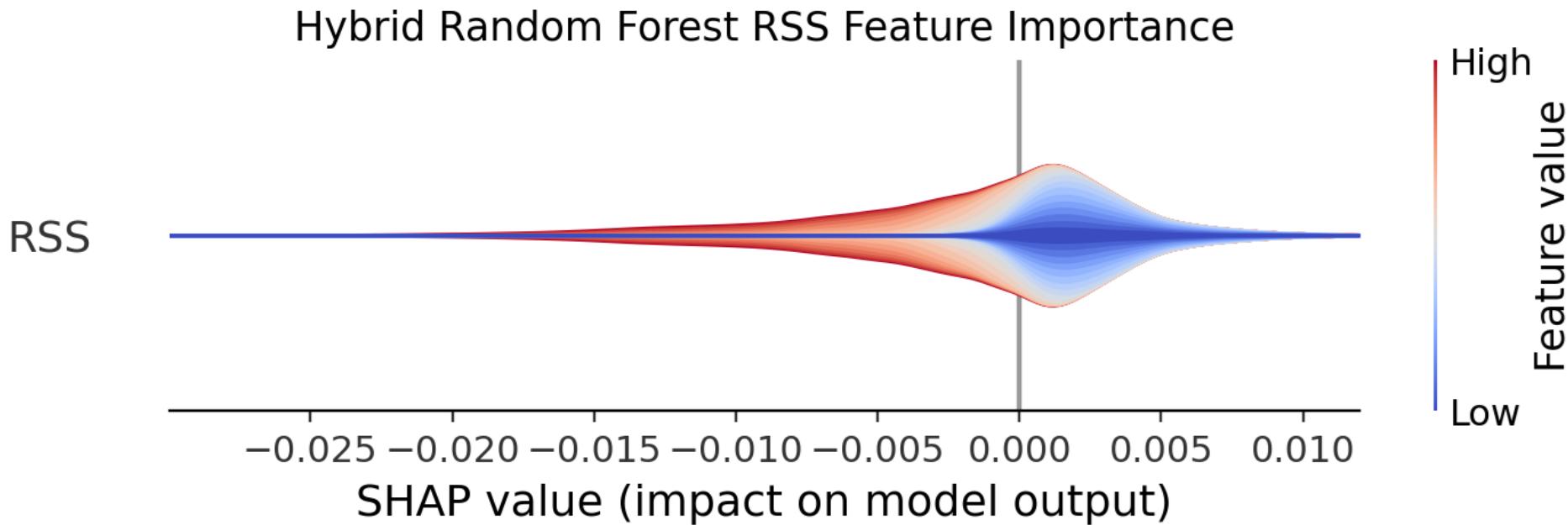
- **Contribution of each feature to the final prediction, averaged across all combinations of other features**
- **Great properties:**
 - Same units as output (probabilities in our case)
 - Additive

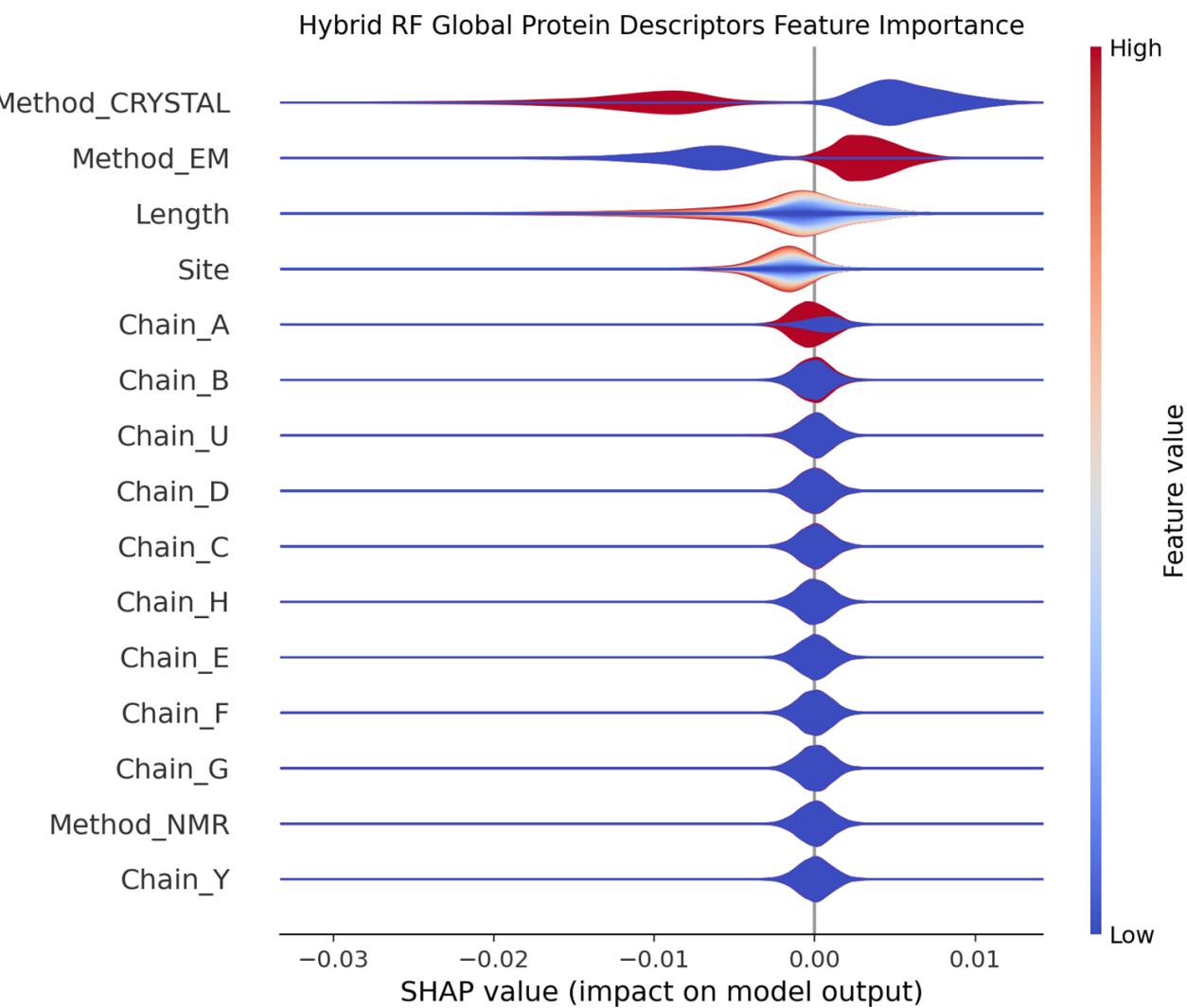
Row-Wise Mean Absolute SHAP By Categories



SHAP Layered Violin Plot

- High feature value + positive SHAP = positive contribution
- High feature value + negative SHAP = negative contribution





Feature Importance Summary

- Functionally important sites tend to be more cryptic
- Method influences crypticity observation
- Residue-level properties are poor predictors

Contents

1. Introduction and background information
2. Dataset
3. Model Training
 - a. Transformer Encoder
 - b. Fusion Encoder
 - c. Single Models
 - d. Hybrid Models
4. Feature Importance
- 5. Outlook**

Limitations and Risks

- Class imbalance between crypticity target variable
- Use of experimental RSA samples
- Poor PTM coverage
- Lack of relevant training features
 - Ligand effects, 3D structure, structural dynamics
- Transformer largely failed to recognize global structural features (monotonous attention maps)
- Chose 25 and 20 for RSA somewhat arbitrarily

Outlook

- Integrate with pretrained structural predictors (e.g. AlphaFold)
- Integrate with molecular dynamics simulations
- Integrate with ligand docking predictors
- Include more comprehensive PTM data
- Include protein functional annotations

References

1. Oleinikovas, Vladimiras; Saladino, G.; Cossins, B. P.; Francesco Luigi Gervasio. Understanding Cryptic Pocket Formation in Protein Targets by Enhanced Sampling Simulations. *JACS* **2016**, *138* (43), 14257–14263. <https://doi.org/10.1021/jacs.6b05425>.
2. Bemelmans, M. P.; Cournia, Z.; Damm-Ganamet, K. L.; Gervasio, F. L.; Pande, V. Computational Advances in Discovering Cryptic Pockets for Drug Discovery. *Curr. Opin. Struct. Biol.* **2025**, *90*, 102975. <https://doi.org/10.1016/j.sbi.2024.102975>.
3. Cimermancic, P.; Weinkam, P.; Rettenmaier, T. J.; Bichmann, L.; Keedy, D. A.; Woldeyes, R. A.; Schneidman-Duhovny, D.; Demerdash, O. N.; Mitchell, J. C.; Wells, J. A.; Fraser, J. S.; Sali, A. CryptoSite: Expanding the Druggable Proteome by Characterization and Prediction of Cryptic Binding Sites. *J. Mol. Bio.* **2016**, *428* (4), 709–719. <https://doi.org/10.1016/j.jmb.2016.01.029>.
4. Meller, A.; Ward, M.; Borowsky, J.; Kshirsagar, M.; Lotthammer, J. M.; Oviedo, F.; Ferres, J. L.; Bowman, G. R. Predicting Locations of Cryptic Pockets from Single Protein Structures Using the PocketMiner Graph Neural Network. *Nat. Commun.* **2023**, *14* (1). <https://doi.org/10.1038/s41467-023-36699-3>.
5. Berman, H.M.; Westbrook, J.; Feng, Z.; Gilliland, G.; Bhat, T.N.; Weissig, H.; Shindyalov, I.N.; Bourne, P.E. The Protein Data Bank. *Nucleic Acids Res.* **2000**, *28*: 235-242 <https://doi.org/10.1093/nar/28.1.235>.
6. Bateman, A.; Martin, M.-J.; Orchard, S.; Magrane, M.; Adesina, A.; Ahmad, S.; Bowler-Barnett, E. H.; Hema Bye-A-Jee; Carpentier, D.; Denny, P.; Fan, J.; Garmiri, P.; Jose, L.; Hussein, A.; Alexandr Ignatchenko; Insana, G.; Rizwan Ishtiaq; Joshi, V.; Dushyanth Jyothi; Swaathi Kandasamy. UniProt: The Universal Protein Knowledgebase in 2025. *Nucleic Acids Res.* **2024**, *53* (D1). <https://doi.org/10.1093/nar/gkae1010>.
7. Chang, K. T.; Guo, J.; di Ronza, A.; Sardiello, M. Aminode: Identification of Evolutionary Constraints in the Human Proteome. *Sci. Rep.* **2018**, *8* (1). <https://doi.org/10.1038/s41598-018-19744-w>.
8. Hornbeck, P. V.; Zhang, B.; Murray, B.; Kornhauser, J. M.; Latham, V.; Skrzypek, E. PhosphoSitePlus, 2014: Mutations, PTMs and Recalibrations. *Nucleic Acids Res.* **2014**, *43* (D1), D512–D520. <https://doi.org/10.1093/nar/gku1267>.
9. Gu, Y.; Tinn, R.; Cheng, H.; Lucas, M.; Usuyama, N.; Liu, X.; Naumann, T.; Gao, J.; Poon, H. Domain-Specific Language Model Pretraining for Biomedical Natural Language Processing. *ACM Trans. Comput. Healthcare* **2022**, *3* (1), 1–23. <https://doi.org/10.1145/3458754>.

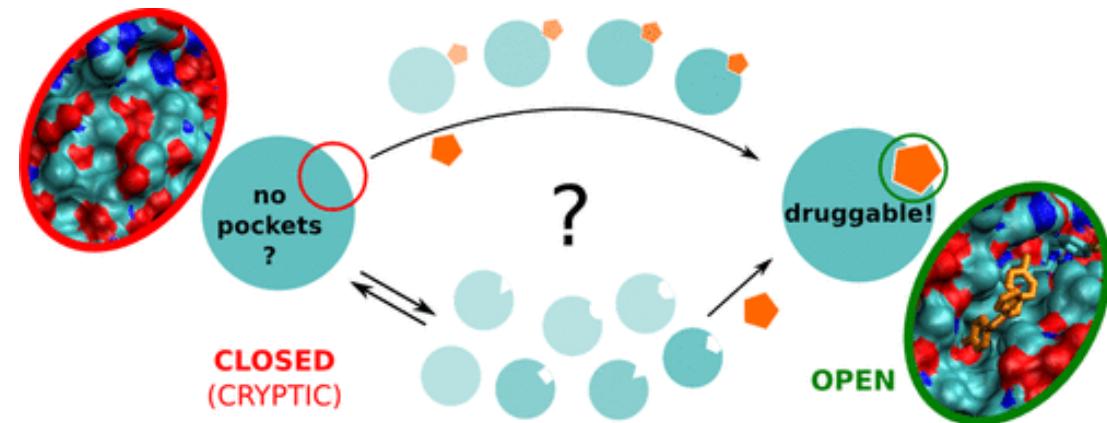
Acknowledgement

Nano Banana (<https://gemini.google/overview/image-generation>) was used for the generation of certain images as noted on the relevant slides.

Supplementary Information

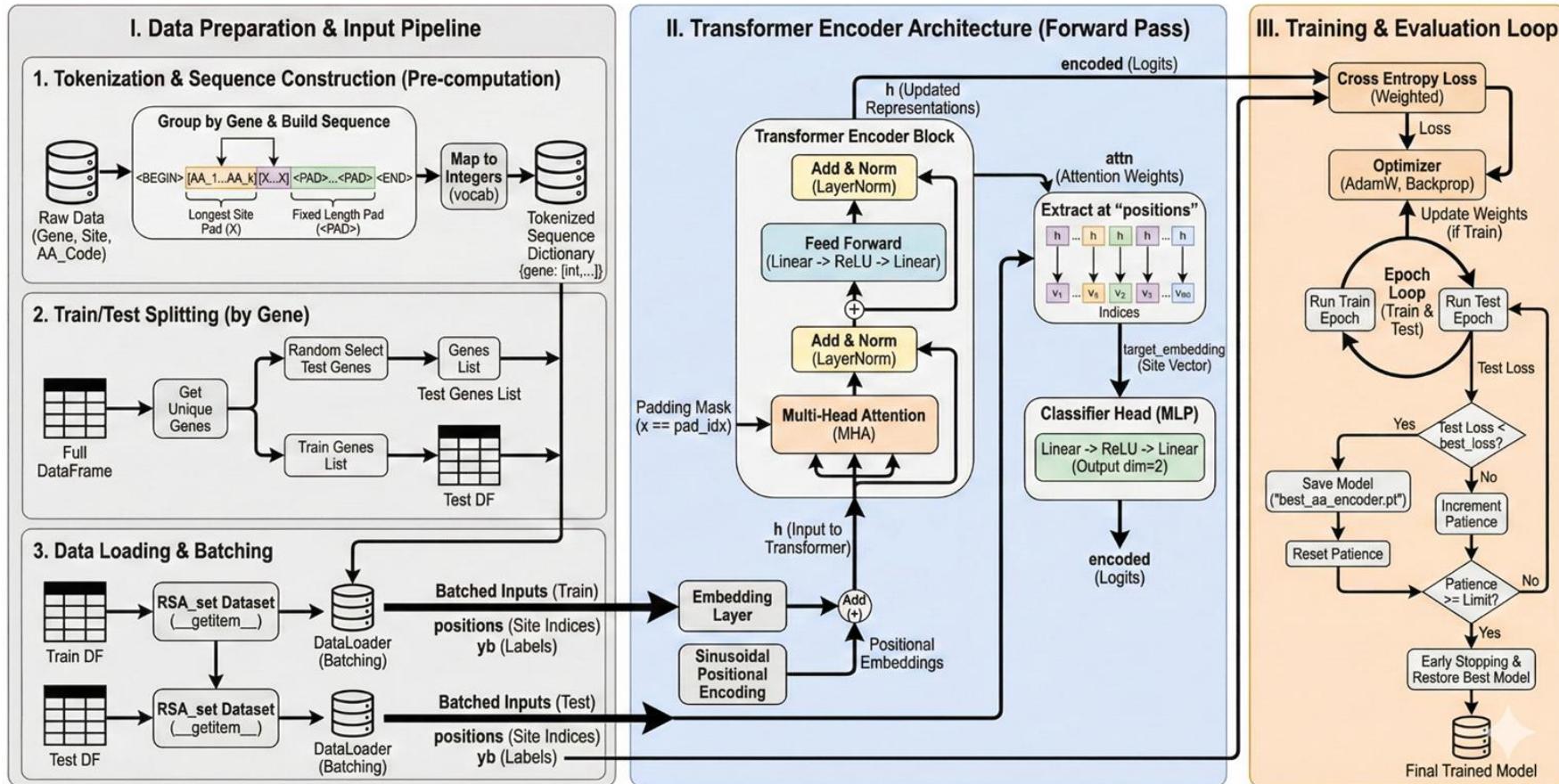
Broad Mechanisms for Cryptic Pocket Opening

- Induced fit
- Conformation selection



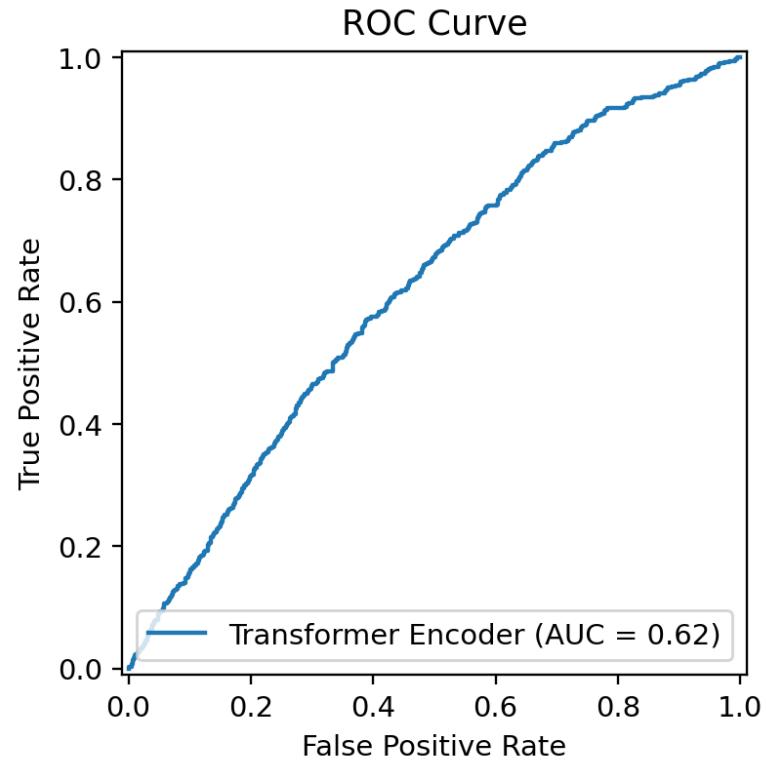
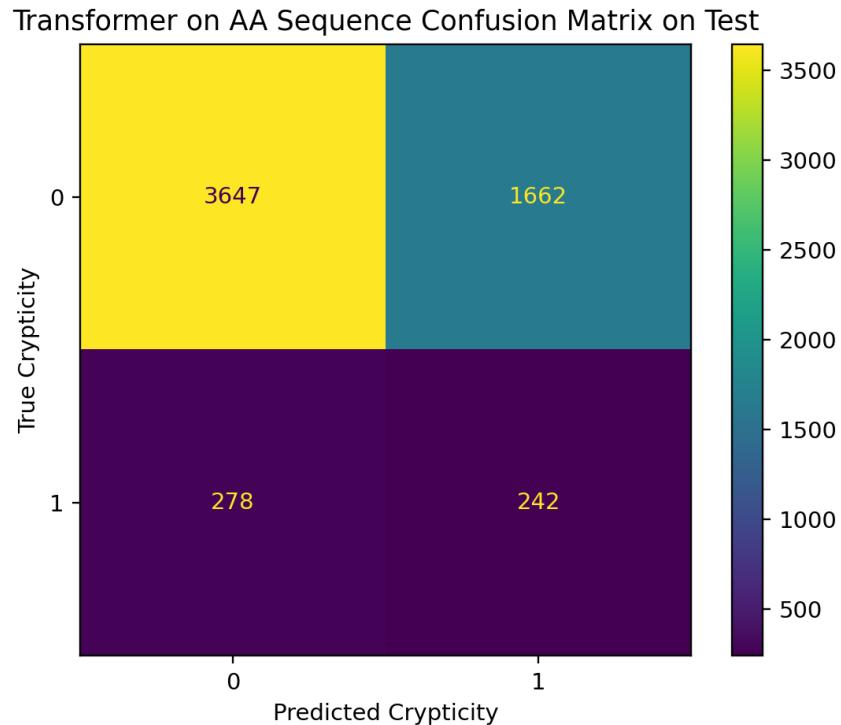
Adapted from Oleinikovas et al.

Transformer Encoder for Protein Cryptic Site Identification: Architecture & Training Flow

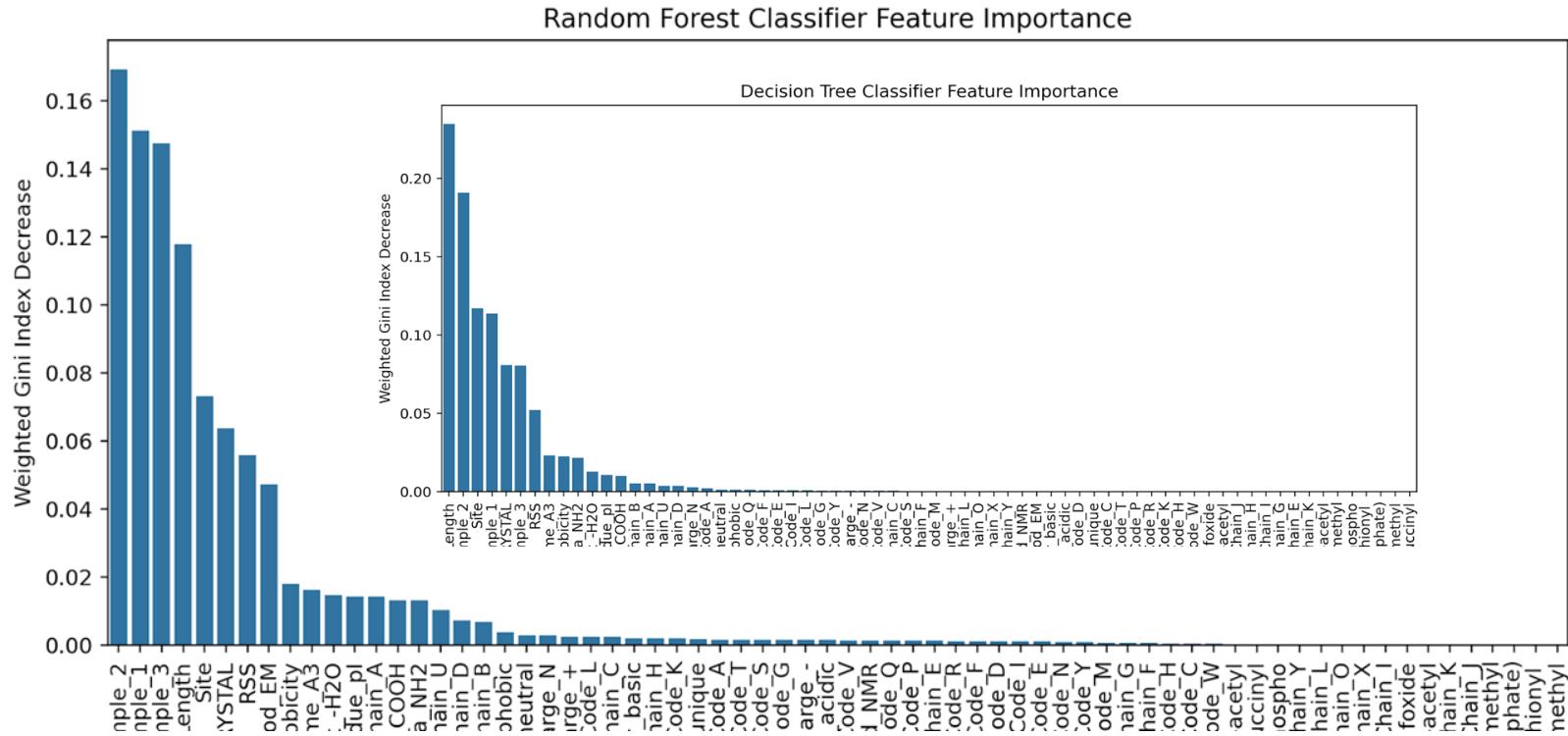


- Accuracy: 0.67
- F1 Score: 0.20
- ROC AUC: 0.618

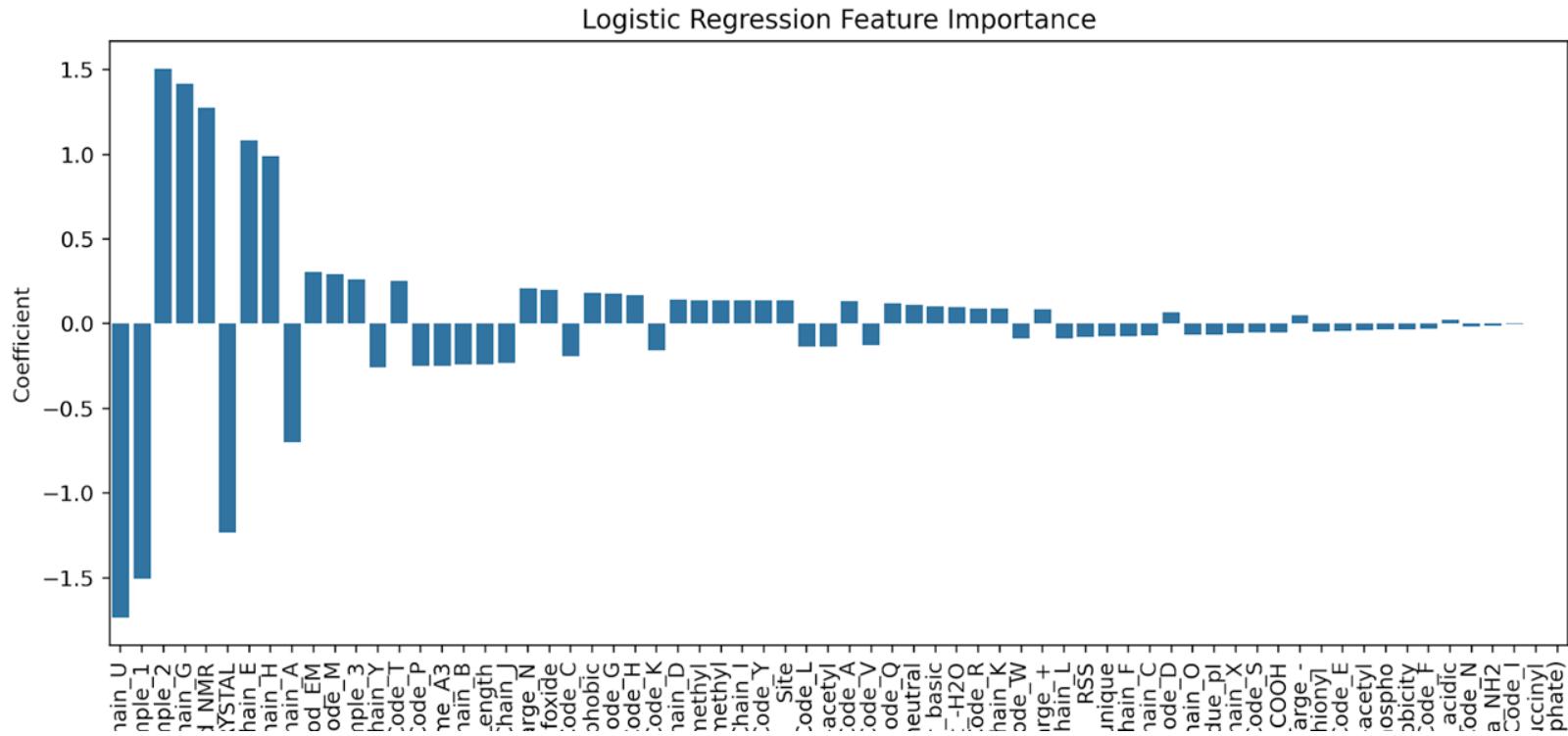
Transformer Encoder Performance



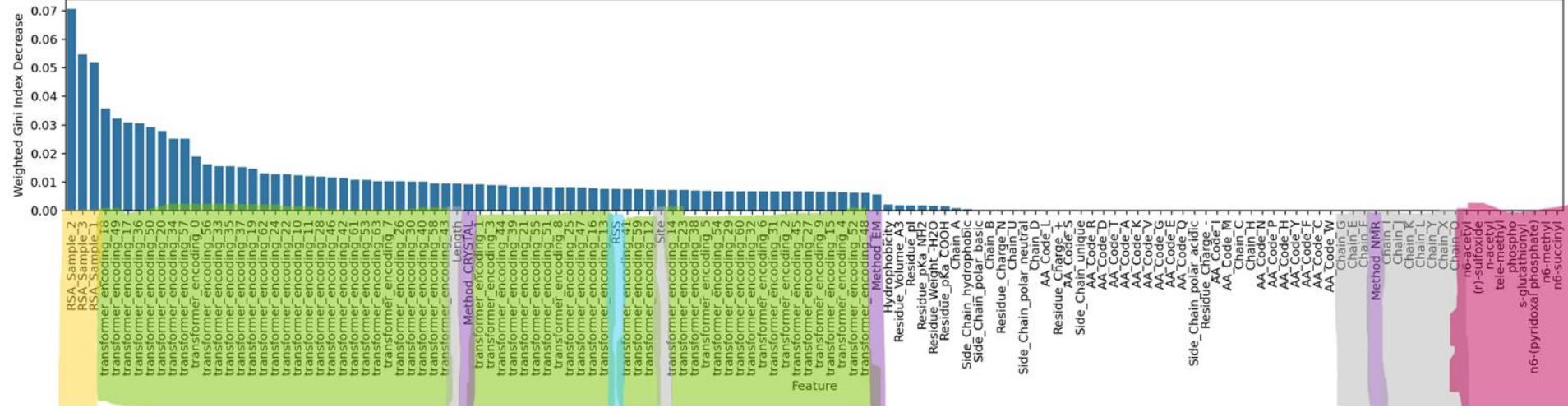
Model dependent: Gini importance



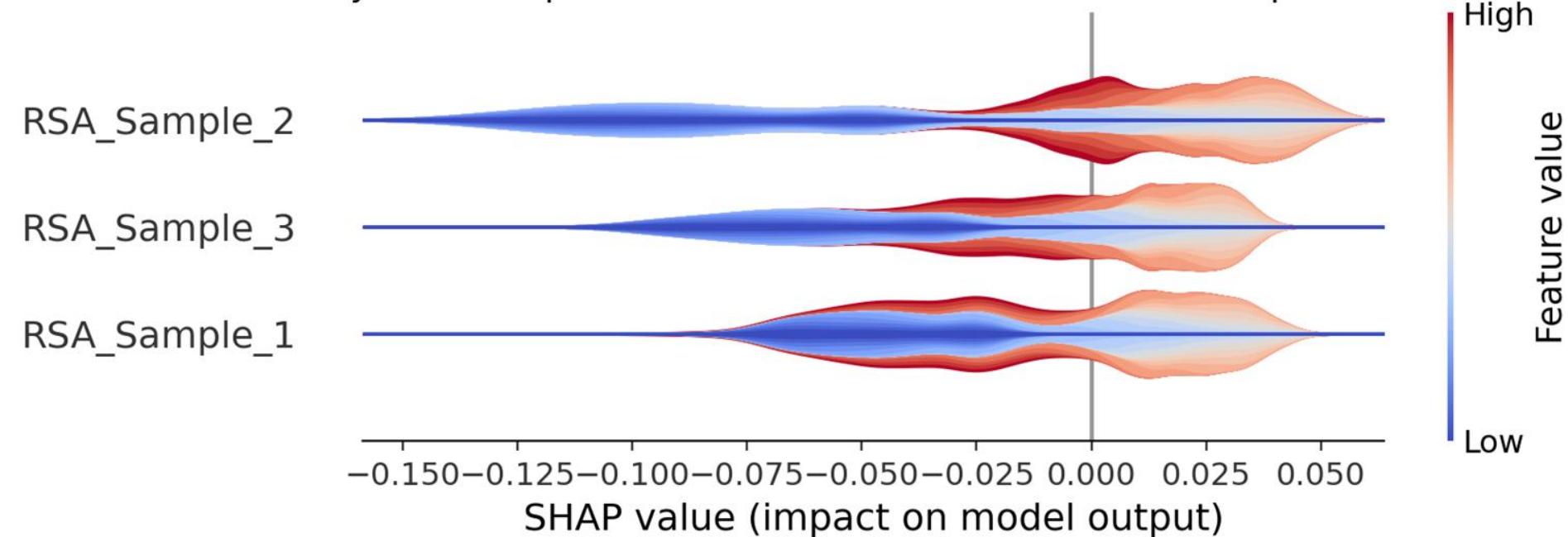
Model dependent: Coefficients



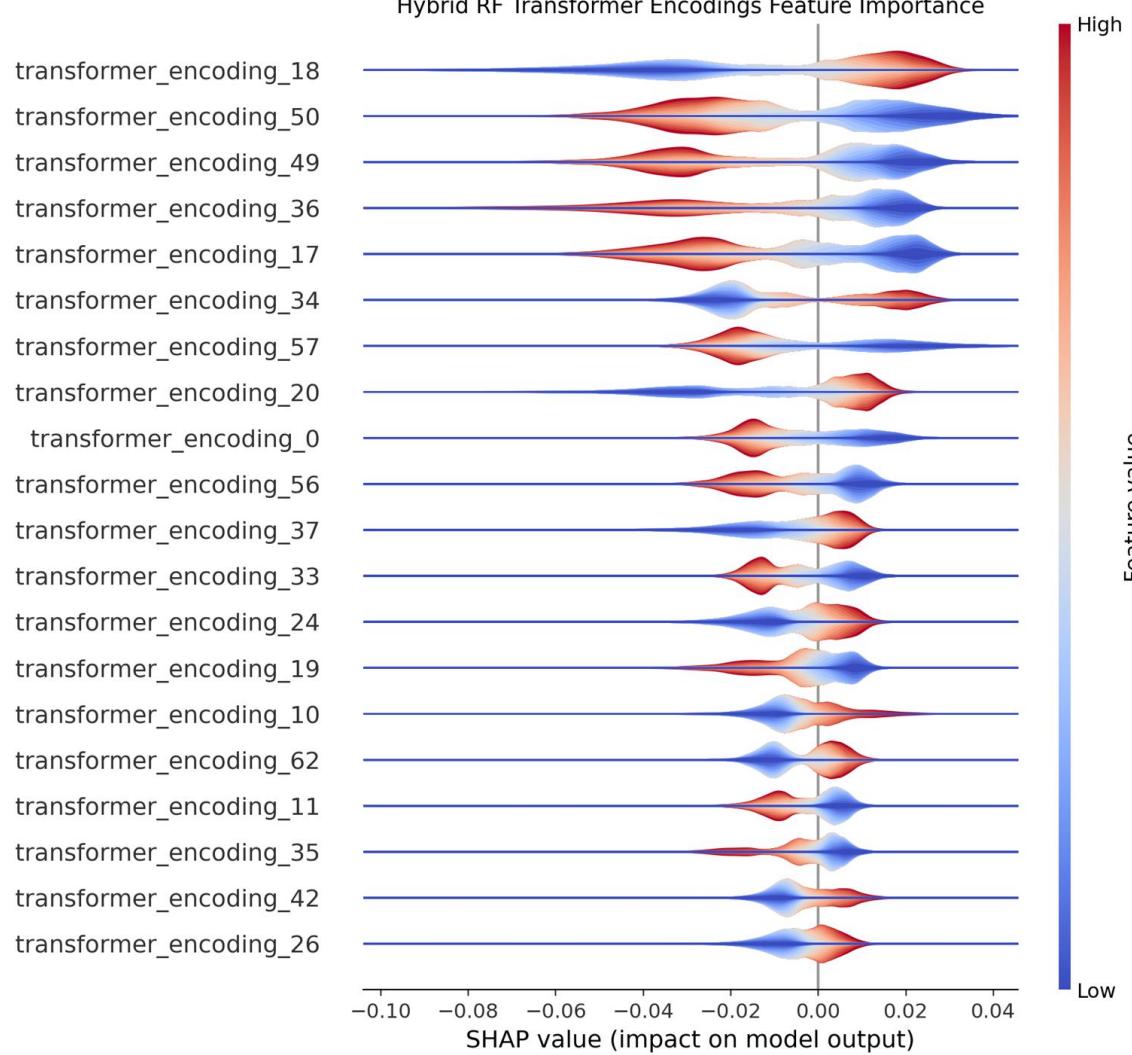
Hybrid Random Forest Classifier Feature Importance

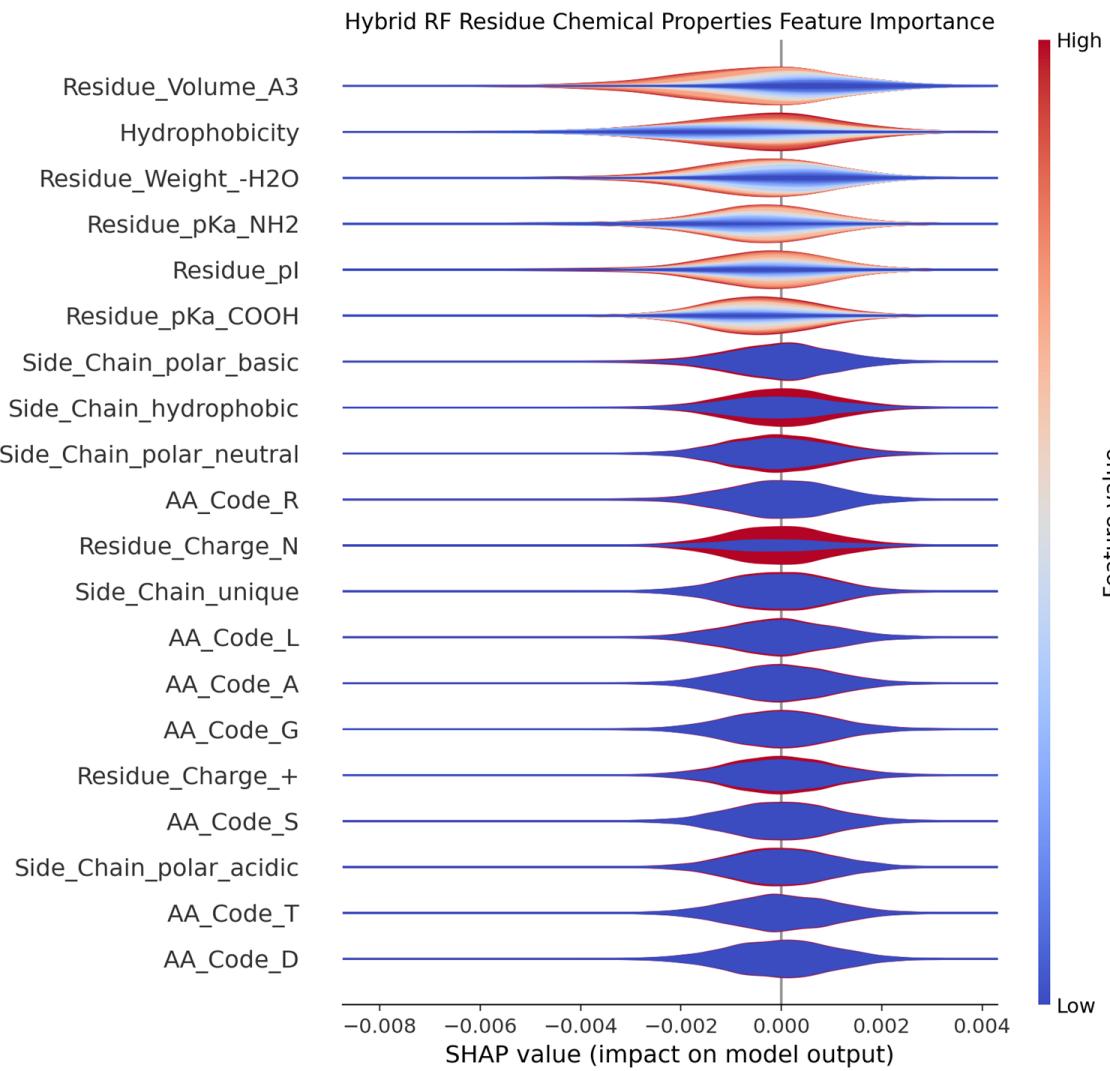


Hybrid RF Experimental RSA Measurements Feature Importance



Hybrid RF Transformer Encodings Feature Importance





Hybrid RF Residue Modification Feature Importance

