

Computers in Biology and Medicine

DHUpredET: A Comparative Computational Approach for Identification of Dihydrouridine Modification Sites in RNA Sequence

--Manuscript Draft--

Manuscript Number:	CIBM-D-24-11894
Article Type:	Full Length Article
Keywords:	Dihydrouridine; Extra tree classifier; RNA molecules; Position-specific two nucleotides; Feature extraction
Corresponding Author:	Md. Shazzad Hossain Shaon Oakland University Michigan State, USA, UNITED STATES
First Author:	Md. Fahim Sultan
Order of Authors:	Md. Fahim Sultan Tasmin Karim Md. Shazzad Hossain Shaon Sayed Mehedi Azim Md. Yeasin Biplob Md. Shoaib Hossain Alshan Iman Dehzangi Mst Shapna Akter
Abstract:	Dihydrouridine (DHU/D) is a modified nucleoside found in RNA molecules. It is formed by reducing the uridine nucleoside, a process that eliminates two hydrogen atoms. DHU has been associated with a variety of diseases, including cancer, brain tumors, and hormonal disorders. Predicting DHU locations in illness-associated RNAs can help us understand disease processes and identify possible treatment targets. However, laboratory-based detection of D sites is laborious and expensive. In this study, we conducted a comparative computational analysis to identify D sites, employing effective machine learning models and efficient feature encoding methods. Initially, we explored various state-of-the-art feature encoding approaches, evaluating 30 machine learning techniques for each, and favored the top eight models based on their independent testing and cross-validation outcomes. As a result, we introduced the DHUpredET, using the extra tree classifier methods for predicting DHU sites. The DHUpredET model demonstrated balanced performance across all evaluation criteria, outperforming state-of-the-art models by 8% and 14% in terms of accuracy and sensitivity on an independent test set. Further analysis revealed that the model achieved higher accuracy with position-specific two nucleotides (PS2) features, leading us to conclude that PS2 features are best suited for the DHUpredET model. Therefore, our proposed model emerges as the most favorite choice for predicting D sites. In addition, we conducted an in-depth analysis of local features and identified a particularly significant attribute with a feature score of 0.035 for PS2_299 attributes. This tool holds immense promise as an advantageous instrument for accelerating the discovery of D modification sites, which contributes to many targeting therapeutic and understanding RNA structure. The datasets used in the study and the source codes are publicly available at https://github.com/Shazzad-Shaon3404/DHUpredET-DHU-prediction.git .
Suggested Reviewers:	
Opposed Reviewers:	

Declaration of interests

- The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.
- The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

none

30 September 2024

Dear,

Editor-in-Chief

Computers in Biology and Medicine, Elsevier.

I am writing to submit my manuscript entitled "**DHUpredET: A Comparative Computational Approach for Identification of Dihydrouridine Modification Sites in RNA Sequence.**" for consideration for publication in the journal "**Computers in Biology and Medicine**".

In this paper, we present an additional and effective model for detecting Dihydrouridine Modification Sites (D/ DHU) based on nine feature encoding techniques. Our model, called DHUpredET, outperformed existing approaches in terms of accuracy and effectiveness in discovering ACPs. Ultimately, this model aids in detecting DHU site more quickly and precisely, saving time and expense.

We believe that our work will be of interest to the readers of "Computers in Biology and Medicine" and will make a valuable contribution to the field of cancer treatment.

Type of Manuscript: Research article.

*** Correspondence**

Md. Shazzad Hossain Shaon, Department of Computer Science and Informatics, Oakland University, Rochester, MI 48309, USA.

Thank you for considering our submission. We look forward to hearing from you.

Sincerely,

Md Shazzad Hossain Shaon

Department of Computer Science and Informatics

Email: shazzad15-3404@diu.edu.bd

Highlight

1. Employed nine feature encoding approaches: modern NLP features and biologically significant properties.
2. Developed DHUpredET, a machine learning-based method which outperforms earlier approaches.
3. Built a simpler model and evaluated it using numerous evaluation criteria to ensure robust and consistent performance.
4. Determined which characteristics are most relevant for prediction, with (PS2) appearing as the most critical.

1 **DHUpredET: A Comparative Computational Approach for Identification of**
2 **Dihydrouridine Modification Sites in RNA Sequence**

3 Md. Fahim Sultan ^a, Tasmin Karim ^a, Md. Shazzad Hossain Shaon ^{a*}, Sayed Mehedi Azim ^b, Md.
4 Yeasin Biplob ^c, Md. Shoaib Hossain Alshan ^d Iman Dehzangi ^{b,e}, Mst Shapna Akter ^f

- 5 a. Department of Computer Science and Informatics, Oakland University,
6 Rochester, MI 48309, USA.
7 b. Center for Computational and Integrative Biology, Rutgers University, Camden, NJ 08102,
8 USA.
9 c. Department of Computer Science & Engineering, Daffodil International University,
10 Daffodil Smart City, Birulia, Dhaka, 1216, Bangladesh.
11 d. Department of Computer Science, Rutgers University, Camden, NJ, 08102, USA
12 e. Department of Computer Science and Engineering, Shanto-Mariam University of Creative
13 Technology, Uttara, Dhaka-1230, Bangladesh.
14 f. Department of Computer Science and Engineering, Oakland University, Michigan, USA.

15 **Email in order:**

- 16 fahim15-3416@diu.edu.bd,
17 tasmin15-2920@diu.edu.bd,
18 shazzad15-3404@diu.edu.bd,
19 sayedmehedi.azim@rutgers.edu ,
20 yeasin15-3055@diu.edu.bd ,
21 shoaibhossainalshan@gmail.com
22 i.dehzangi@rutgers.edu ,
23 akter@oakland.edu

24

25

26 **Abstract-** Dihydrouridine (DHU/D) is a modified nucleoside found in RNA molecules. It is
27 formed by reducing the uridine nucleoside, a process that eliminates two hydrogen atoms. DHU
28 has been associated with a variety of diseases, including cancer, brain tumors, and hormonal
29 disorders. Predicting DHU locations in illness-associated RNAs can help us understand disease
30 processes and identify possible treatment targets. However, laboratory-based detection of D sites
31 is laborious and expensive. In this study, we conducted a comparative computational analysis to
32 identify D sites, employing effective machine learning models and efficient feature encoding
33 methods. Initially, we explored various state-of-the-art feature encoding approaches, evaluating
34 machine learning techniques for each, and favored the top eight models based on their independent
35 testing and cross-validation outcomes. As a result, we introduced the DHUpredET, using the extra
36 tree classifier methods for predicting DHU sites. The DHUpredET model demonstrated balanced
37 performance across all evaluation criteria, outperforming state-of-the-art models by 8% and 14%
38 in terms of accuracy and sensitivity on an independent test set. Further analysis revealed that the
39 model achieved higher accuracy with position-specific two nucleotides (PS2) features, leading us
40 to conclude that PS2 features are best suited for the DHUpredET model. Therefore, our proposed
41 model emerges as the most favorite choice for predicting D sites. In addition, we conducted an in-
42 depth analysis of local features and identified a particularly significant attribute with a feature
43 score of 0.035 for PS2_299 attributes. This tool holds immense promise as an advantageous
44 instrument for accelerating the discovery of D modification sites, which contributes to many
45 targeting therapeutic and understanding RNA structure. The datasets used in the study and the
46 source codes are publicly available at <https://github.com/Shazzad-Shaon3404/DHUpredET-DHU-prediction.git>.
47

48

49 **Keywords:** Dihydrouridine, Extra tree classifier, RNA molecules, Position-specific two
50 nucleotides, Feature extraction

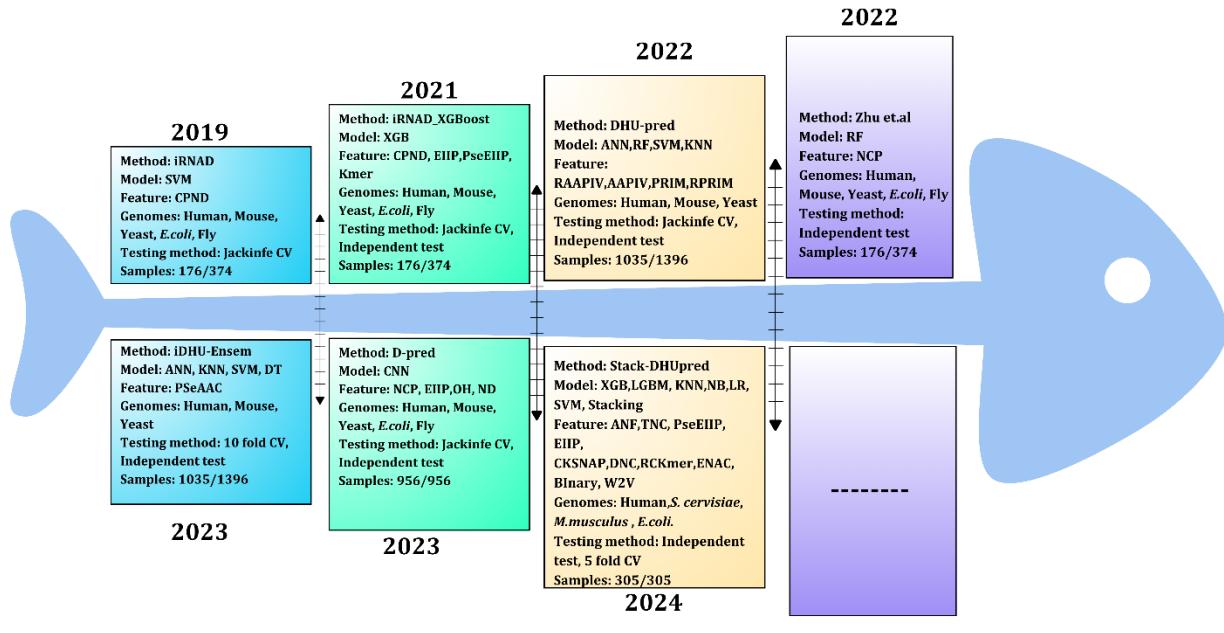
51 **Introduction**

52 RNA Modification refers to a sequence of chemical changes made to initial ribonucleic acid
53 (RNA) transcripts that culminate in the synthesis of mature RNA molecules essential to a variety
54 of biological activities [1]. At present, more than 300 particular types of RNA modifications are
55 discovered [1-5]. Furthermore, the environment of RNA modification is strikingly similar across
56 all three domains of life: eukaryotes, bacteria, and archaea. Despite their great evolutionary
57 deficiencies, these domains contain numerous alterations and modifying enzymes [1, 6, 7].

58 Dihydrouridine (DHU), denoted by the letter "D" in RNA sequences and having the chemical
59 formula $C_9H_{14}N_2O_6$, is a significant alteration identified in numerous forms of RNA, including
60 transfer RNA (tRNA), messenger RNA (mRNA), and small nucleolar RNA (snoRNA) in different
61 organisms. This alteration has sparked enthusiasm because of its prevalent presence and probable
62 functional importance in RNA molecules [8-10]. The dihydrouridine biosynthesis molecules
63 catalyze the reduction of uridine's (U) C5-C6 double bond, leading to the formation of
64 dihydrouridine [6]. This alteration can be observed at the variable position of the anticodon helix

65 in tRNA molecules and poses major health consequences, including the development of lung
 66 cancer, Alzheimer's disease, and Huntington's disorders [11-13]. Therefore, D sites are crucial for
 67 understanding RNA molecules' structure, function, and regulatory activities, as well as their
 68 prospective uses in diagnostics and therapies [14].

69 Experimental methods for recognizing D sites, such as mass spectrometry or antibody-based
 70 assays, can be complicated, time-consuming, and expensive. On the other hand, computational
 71 approaches provide a high-throughput capabilities alternative, facilitating experts to evaluate
 72 massive RNA sequence databases rapidly and effectively. Advanced computational models,
 73 especially using machine learning and artificial intelligence, can achieve high levels of accuracy
 74 in predicting D sites. They can detect subtle patterns and features in RNA sequences that might be
 75 missed by traditional methods. **Figure 1** presents the fishbone diagram of the existing machine
 76 learning methods proposed to predict D sites.



77
 78 **Figure 1.** An overview of the existing methods in various years, represented with a fishbone
 79 diagram.
 80

81 During the past few years, several machine learning-based methods proposed to predict D sites. In
 82 2019, Xu et al. proposed an iRNAD framework based on a support vector machine (SVM) model
 83 with various features and achieved promising results [14]. Later on, Dou et al. developed
 84 iRNAD_XGBOOST with extreme gradient boosting (XGB) on different feature selection
 85 approaches [15]. In 2022, Zhu et al. introduced a Random Forest (RF) based machine learning
 86 approach for the D sites identification [16]. At the same time, Suleman et al. proposed the DHU-
 87 pred method with RF, where the authors used a relatively larger dataset compared to their prior
 88 studies to train their model [17]. Later on, the iDHU-Ensem method was proposed, where Suleman
 89 et al. included artificial neural network (ANN), K-nearest neighbor (KNN), SVM, and decision
 90 tree (DT), and obtained a better result [18]. In the same time, Yu et al. developed a D-pred model
 91 with a local self-attention layer and a convolutional neural network (CNN) [9]. Self-attention
 92 techniques, particularly local self-attention, might prove computationally expensive, especially

when dealing with complex sequences. Most recently, Roshid et al. developed stacking-based machine learning models called Stack-DHUpred, where the authors used probabilistic values from the baseline models and obtained a better outcome compared to their previous studies [19]. However, there is still room for improvement in the performance.

In general, most of the previous studies used a limited number of feature embedding approaches and utilized small datasets to train their model. As a result, the present study focuses on recently obtained datasets using various feature embedding approaches to determine the best feature extraction method and framework for predicting D sites.

In this study, we executed an in-depth analysis employing various machine learning techniques. We also use nine feature encoding methods, including state-of-the-art methods generated from natural language features, biologically feasible physical and chemical properties, nucleic acid-based compositional features, and residue composition data. Through this extensive analysis, we sought to uncover the most informative features and optimal machine-learning algorithms for predicting D sites in RNA sequences based on various types of evaluation metrics. As a result, we propose DHUpredET, a machine learning-based method than can significantly outperform previous studies found in the literature. Notably, it excels at leveraging position-specific two nucleotides (PS2) features with appropriate execution. DHUpredET identifies the positive class with greater than 85% and specifies the negative class with more than 82% accuracy, which produces a balance outcome. The datasets used in the study and the source code for DHUpredET are publicly available at <https://github.com/Shazzad-Shaon3404/DHUpredET-DHU-prediction.git>.

2. Materials and methods

2.1. Dataset descriptions

In this study, we use the dataset that was introduced in [19]. This dataset contains a total of 805 active compounds from different species, such as *Homo sapiens*, *Mus musculus*, *Saccharomyces cerevisiae*, and *Escherichia coli* [3,20,10]. Next, they used a 90% cluster database at high identity with tolerance (CD-HIT) method [21] to reduce the redundant data of the positive molecules. As a result, the number of positive samples was reduced to 305 samples. They also collected same number of negative samples (305) from different studies [22-24]. Finally, they used, 244 positive and 244 negative compounds as the training set, while an additional set of 61 positive and 61 negative samples were reserved as the testing set. **Table 1.** Describe the details of the datasets.

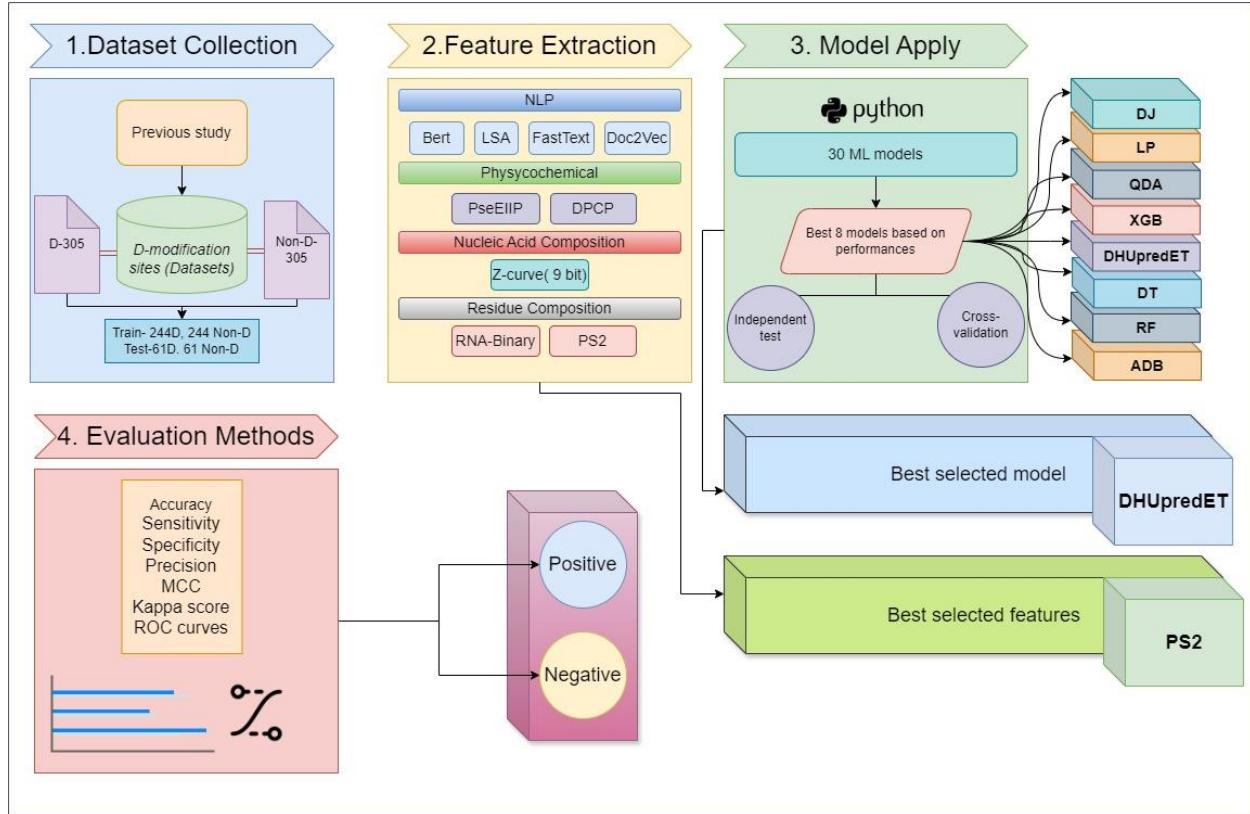
Table 1. Dataset description of the train file and test file used in this study.

Mode	Total	Train	Test
D sites (positive)	305 D sites	244 D sites	61 D sites
non-D sites (negative)	305 non-D sites	244 non-D sites	61 non-D sites

2.2. Overall methodologies
This study presents a machine learning-based approach for DHU prediction. **Figure 2** demonstrates the overall methodology of the study. At first, we collected the dataset through evaluation of previous studies, resulting in an extensive variety of RNA sequences from different species. Then, we used multiple feature embedding techniques on the collected data. In order to improve efficiency, we utilized a variety of machine learning models. As a result, the eight greatest

136 performing models were chosen after a rigorous evaluation of their overall performances. Among
 137 them, DHUpredET emerged as the most promising model for predicting D sites and PS2 as the
 138 most effective encoding methods.

139

140
141

142 **Figure 2.** Workflows of the current study. Dataset collection from the previous studies, feature
 143 extraction approach of four kinds of descriptors with nine feature encoding methods, application
 144 of the various models, and selection process, the overall comparison of the models, and selecting
 145 the optimal feature and best-fit models for the study.

146

147 2.3. Feature extraction

148 Feature extraction is vital step for a better performance of machine learning models because it
 149 transforms raw data into a more appropriate representation that allows for successful learning and
 150 prediction. In this study, we use various feature extraction techniques to determine the most
 151 effective attributes for D sites prediction. These techniques included NLP-based methods such as
 152 bidirectional encoder representations from transformers (BERT) [25-28], FastText word
 153 embedding (FastText) [29,30], latent semantic analysis (LSA) [31-33], and document-to-vector
 154 (Doc2Vec) [34-36]. Furthermore, physicochemical characteristics such as electron-ion interaction
 155 pseudopotentials of trinucleotide (PseEIIP) [37, 38] and dinucleotide physicochemical properties
 156 (DPCP) [39] are also deployed. The study further utilized residue composition-based RNA binary
 157 (Binary) [40, 41] and PS2 [42,43], in addition to nucleic acid composition-based features,
 158 including Z-curve phase-specific mononucleotides (Z-curve 9-bit) [44]. These strategies are
 159 intended to improve the accuracy and efficacy of D site detection. In the following subsections, all
 160 these methods are described, in detail.

161
162
163
164

165 **2.3.1 NLP-based feature extractions**

166
167 **BERT:** Bert is well-known in the data mining sector for making encoder transformers. With a pre-
168 trained approach, it has received a lot of attention and recognition. It encodes contextual by
169 analyzing words bidirectionally, indicating that it considers both left and right circumstances. It
170 tokenizes input text, fine-tunes for certain tasks, and supplies complemented word visualizations,
171 allowing it to excel at various text-based activities [25-28]. The formula of the BERT feature can
172 be expressed as:

173
$$B = \prod_{m \in Ma} P(m | W, A) \quad (1)$$

174
$$P(m | W, A) = X(W \cdot O_m) \quad (2)$$

175 In equations 1 and 2, Ma is the vector of the masked word in the sentence, W is the vector
176 associated with a word present in the input sentence, A is the parameter of the approach, O_m
177 signifies the output embedding for the masked word, m , $W \cdot O_m$ dot product computed between the
178 word vector and the outcome embedding vector, and $P(m | W, A)$ refers to the probability of the
179 word m , conditioned on both W and A .

180
181 **FastText:** FastText is an NLP model developed by Facebook AI Research that excels at word
182 embedding and text classification tasks. FastText describes languages by organizing symbols
183 across n-grams. As an instance, the pattern of "XYZM" can be split into smaller n-grams such as
184 "XY," "XYZ," "XYM," "YZ," "YZM," and "ZM". The n-grams preserve mathematical data which
185 allows them to encode sentences with equivalent frequencies [29, 30]. The equation can be stated
186 as:

187
188

$$F = -\frac{1}{T} \sum_{n=1}^T x_n \log(P(TL_{a_n})) \quad (3)$$

189 In equation 3, F denotes the loss function. It averages a negative logarithm of the anticipated
190 distribution of probabilities P , throughout a sequence of n-gram features denoted as x_n , each of
191 which can be represented by word embeddings of a_n . The look-up matrices of word embeddings
192 are denoted by L . T represents linear output transformation. The SoftMax function's coefficient is
193 applied to the output.

194
195 **LSA:** LSA is the most popular NLP-based embedding system in the data mining fields. It works
196 by initially displaying a collection of text documents as a matrix, with rows representing distinct
197 terms (words) and columns representing documents. Then, it performs singular value
198 decomposition (SVD) on the matrix. LSA decreases the matrix's dimensionality by maintaining
199 just the top k singular values and accompanying singular vectors, resulting in a lower-dimensional
200 semantic space. Each phrase and document is represented as a vector in this semantic space, which
201 captures latent semantic associations [31-33]. The mathematical term of the LSA feature denotes:

202
203

$$LSA = U \sum V_N \quad (4)$$

204 Where, U is the left singular vectors, representing the relationships between terms in the
 205 reduced-dimensional semantic space. V_N is the transpose of matrix V , representing associations
 206 among texts in a reduced-dimensional conceptual environment.

207 **Doc2Vec:** Doc2Vec is a sophisticated tool for creating document embeddings that capture the
 208 semantic content of documents, rendering it suitable for various text analysis purposes. It was
 209 introduced by researchers at Google and is based on the distributed memory (DM) and distributed
 210 bag of words (DBOW) architectures. It operates by learning a neural network to determine the
 211 wider context of a document, whereas word-to-vector analyzes the surrounding words of a
 212 particular word. Each document is allocated a distinctive vector in a space with high dimensions
 213 that captures the semantic meaning [34-36]. The formula can be stated as:

$$214 \quad D = N1 \sum i = 1N(a + V) \quad (5)$$

215 Here, D is the main vector, N is the number of words, a represents the vectors of each word in the
 216 document, V refers to the particular vector, and $N1$ is the pooling operation.

217 2.3.2 Physicochemical-based feature extractions

218 **DPCP:** This is the structural and chemical properties of neighboring nucleotide pairs inside RNA
 219 sequences. This approach sheds light on RNA's stability, flexibility, and functional qualities by
 220 considering different aspects, including base stacking, hydrogen bonding, and nucleotide
 221 interactions. [39]. The formula is referred to as:

$$222 \quad D_{ij} = f(P_i, Q_j) \quad (6)$$

223 where D_{ij} is the combined property of nucleotides, P_i represents the physicochemical values of the
 224 i^{th} nucleotide base in the dinucleotide pair, Q_j is the property value of the j^{th} base, and f is a function
 225 that combines the properties of two bases for the overall property.

226 **PseEIIP:** This feature introduces attributes of modifications to the original EIIP values to improve
 227 their efficacy in identifying nucleotide sequences. The actual formula for determining PseEIIP
 228 values varies based on the individual alterations and characteristics used in the computation.
 229 PseEIIP values are often calculated by combining EIIP values ($A = 1.1260$, $T/U=0.1335$,
 230 $G=0.0806$, $C=0.1340$). These values are based on each base's electron-ion interaction potential.
 231 EIIP values are useful for a variety of statistical inquiries, such as sequence alignment, motif
 232 finding, and identifying functional components in nucleic acid sequences [37-38]. The formula is
 233 represented as:

$$234 \quad EIIP = \frac{\sum_{i=1}^n E_i \cdot N_i}{\sum_{i=1}^n N_i} \quad (7)$$

$$235 \quad PseEIIP = \sum_{i=1}^n (I_i \cdot E_i) \quad (8)$$

236 In equation 7, n is the total number of nucleotides, E_i is the electron-ion values of the i^{th} position,
 237 and N_i is the frequency of the i^{th} .

238 In equation 8, I_i represents weights assigned to each nucleotide base or feature in the sequence, E_i
 239 is the electron-ion values of the i^{th} position, and n is the total number of nucleotides.

240

241

242

243 **2.3.3 Residue composition-based feature extractions**

244 **Binary Profile:** For this feature, each amino acid in the sequence is assigned to a 20-dimensional
 245 array with just one non-zero component. This mapping generates a matching matrix, where each
 246 row represents an amino acid and each column represents a specific property [40, 41].

247 **PS2:** The PS2 matrix is commonly 16×16 , with each row and column representing one of the
 248 available dinucleotides (16 in total: AA, AC, AG, AU, CA, CC, CG, CU, GA, GC, GG, GU, UA,
 249 UC, UG, and UU). The value in each cell of the matrix represents the frequency or likelihood of
 250 the associated dinucleotide pair occurring at that location in the RNA sequence [42, 43]. The
 251 equation can be represented as:

$$252 M_{ij} = \frac{\text{count}(d_{ij})}{N-1} \quad (9)$$

253 Where, M_{ij} represents the frequency, $N-1$ is the total number of positions, and $\text{count}(d_{ij})$ is the
 254 number of dinucleotide pairs of d_{ij} at position i .

255 **2.3.4 Nucleic acid composition-based feature extractions**

256 **Z-curve 9-bit:** Each nucleotide base (A = 100000000, C = 010000000, G = 001000000, and U =
 257 000100000) in the Z-curve 9-bit representation is translated to a 9-bit binary vector that encodes
 258 features like local stacking energy, hydrogen bonding, and base-pairing preferences. These
 259 features are determined using the positional connections of neighboring nucleotides [44].

260 **2.4 Machine Learning Models Construction**

261 In this study, different machine-learning models are used to predict D sites. After extensive testing
 262 with different hyperparameters, eight models stood out for their outstanding performance. These
 263 include AdaBoost (ADB), label propagation (LP), quadratic discriminant analysis (QDA), extreme
 264 gradient boosting (XGB), Decision Tree (DT), Random Forest (RF), Decision Jungle (DJ), and an
 265 extra tree classifier-based approach (ET), referred to as DHUpredET for PS2 feature extraction
 266 procedure as it performed better[45-54]. These models were chosen based on their ability to
 267 reliably predict D sites, assessed by an exhaustive analysis and juxtaposition of their effectiveness
 268 indicators. In **Table 2**, we illustrate all the hyperparameters used for these classifiers.

269 **Table 2.** Hyperparameters description of the study

Models	Best parameters
ADB	estimator=None, *, n_estimators=50, learning_rate=1.0, algorithm='SAMME.R', random_state=None

LP	kernel='knn',max_iter=1000
QDA	*, priors=None, reg_param=0.0, store_c covariance=False, tol=0.0001
XGB	n_estimators=100, max_depth=100, learning_rate=0.1, subsample=1.0, colsample_bytree=1.0, reg_alpha=30, reg_lambda=30, gamma=0, min_child_weight=1
DT	max_depth=10, criterion='entropy', min_samples_split=10, min_samples_leaf=1, max_features=15, random_state=42
RF	n_estimators=100,criterion='entropy',max_features="sqrt",random_state=100
DHUpredET	n_estimators=100, random_state=10, max_depth=None, min_samples_split=4, min_samples_leaf=1, max_features='auto', bootstrap=False, class_weight=None, criterion='gini', n_jobs=None

270

271 **2.4.1 DHUpredET model development procedure**

272

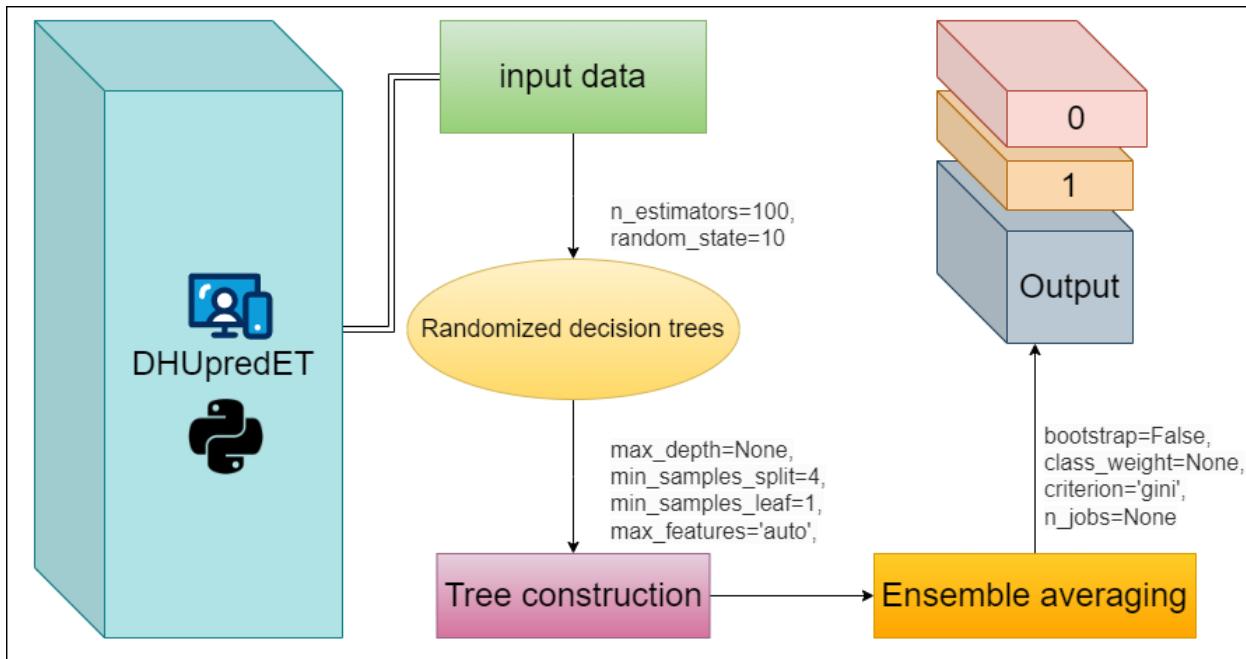
273 This study used a wide range of algorithms, including ten classical approaches, ten voting-based
 274 meta-learning methods, and ten stack-based meta-learning strategies, which are shared in
 275 supplementary materials and provided in the GitHub repository. We chose eight models for their
 276 reasonable scores over numerous assessment measures. It was observed that while other models
 277 exhibited some improvement in accuracy, they displayed a tendency towards bias in positive class
 278 identification. As a result, we selected the Extra Trees Classifier which outperformed all other
 279 classifiers using PS2 features.

280

281 The extra tree classifier is a decision tree-based ensemble learning approach that integrates many
 282 trees to increase forecast accuracy and resilience. In this design, the model is configured to employ
 283 100 decision trees (n_estimators=100), which ensures a varied range of classifiers and reduces
 284 overfitting. The random_state option is set to 10, guaranteeing that the findings are reproducible
 285 over several runs. With max_depth=None, the decision trees can grow until all leaves are pure or
 286 there are less than four samples (min_samples_split=4) at a node, with at least one sample required
 287 at each leaf (min_samples_leaf=1). The max_features option is set to 'auto', indicating that all
 288 features are evaluated for determining the optimal split at each node. Bootstrap sampling is
 289 deactivated (bootstrap=False), indicating that the full dataset is used to build each tree. The
 290 class_weight contention is set to None, which means that all classes are considered equally. The
 291 criterion option is set to 'gini', which specifies the criteria for dividing nodes. Finally, setting
 292 n_jobs=None indicates that the model will use all available processors for parallel processing.

293 **Figure 3.** Demonstrates the overall procedure of the DHUpredET model.

294



295
296
297
298
299
300
301
302
303
304
305
306

Figure 3. Overall strategies of the DHUpredET model.

2.5 Evaluation Metrics

The study used various evaluation metrics to properly evaluate the models' performance. The metrics used were the Matthews Correlation Coefficient (MCC), sensitivity (Sen), specificity (Spe), precision (Pre), recall (Rec), F1 score (F1s), kappa score (Kpp), and accuracy (Acc) [56-60]. These metrics are formulated as follows:

$$\left\{
 \begin{aligned}
 Acc &= \frac{(TP+TN)}{(TP+TN+FP+FN)} \\
 Spe &= \frac{TN}{(TN+FP)} \\
 Sen &= \frac{TP}{(TP+FN)} \\
 MCC &= \frac{TP*TN-FP*FN}{\sqrt{TP+FP*(TP+FN)*(TN+FP)}} \\
 Pre &= \frac{TP}{(TP+FP)} \\
 Rec &= \frac{TP}{(TP+FN)} \\
 F1s &= \frac{2*(Pre*Rec)}{((Pre+Rec))} \\
 Kpp &= \frac{2*(TP*TN-FP*FN)}{(TP+FP)*(FP+TN)+(TP+FN)*(FN+TN)}
 \end{aligned} \right. \quad (10)$$

307
308
309
310

Where TP represents true positive, TN represents true negative, FP represents false positive, and FN represents false negative.

311 **3. Experimental results**

312 In order to identify the best model for D site prediction, we conducted a thorough analysis,
 313 evaluating a wide range of models. Our goal was to identify the best model and features suitable
 314 for the prediction of D site. **Tables 3** provides a complete performance evaluation of our models.
 315 The corresponding results for 5-fold cross-validation is provided as the supplementary material in
 316 **Table_S1**.

317

318 **Table 3.** Independent test performances of eight models

Descriptor	Classifier	Acc (present in percent age with 1 decimal point)	MCC (Decimal number with 2 decimal point prediction)	Kpp (Decimal number with 2 decimal point prediction)	Pre (present in percentage with 1 decimal point)	Rec (present in percentage with 1 decimal point)	F1s (Decimal number with 2 decimal point prediction)	Sen (present in percentage with 1 decimal point)	Spe (present in percentage with 1 decimal point)
DPCP	RF	65.6	0.31	0.31	66.1	63.9	0.65	67.2	0 63.9
	ADB	69.6	0.39	0.39	66.1	77.0	0.71	62.3	77.0
	DT	54.9	0.09	0.09	54.5	50.9	0.56	50.8	59.0
	XGB	66.3	0.32	0.32	65.6	68.8	0.67	63.9	68.8
	QDA	59.8	0.19	0.19	61.1	54.1	0.57	65.5	54.1
	DJ	66.3	0.32	0.32	65.1	70.4	0.67	62.3	70.4
	LP	66.3	0.32	0.32	66.6	65.5	0.66	67.2	65.5
	ET	69.7	0.39	0.39	70.0	68.8	0.69	70.4	68.8
Z curve 9 bit	RF	65.5	0.31	0.31	63.7	72.1	0.67	59.0	72.1
	ADB	63.1	0.26	0.26	62.5	65.5	0.64	60.6	65.5
	DT	51.6	0.03	0.03	51.7	49.1	0.50	54.1	49.1
	XGB	59.0	0.19	0.18	65.7	37.7	0.47	80.3	37.7
	QDA	65.5	0.31	0.31	64.6	68.8	0.66	62.3	68.8
	DJ	59.8	0.20	0.19	62.0	50.8	0.55	68.8	50.8
	LP	63.1	0.26	0.26	62.1	67.2	0.64	59.0	67.2
	ET	65.5	0.31	0.31	65.0	67.2	0.66	63.9	67.2
Binary	RF	72.1	0.45	0.44	77.5	62.3	0.69	81.9	62.3
	ADB	72.1	0.44	0.44	71.4	73.7	0.72	70.4	73.7
	DT	68.8	0.37	0.37	70.9	63.9	0.67	73.7	63.9
	XGB	72.9	0.49	0.45	86.8	54.1	0.66	91.8	54.1
	QDA	69.6	0.48	0.39	96.1	40.9	0.57	98.3	40.9
	DJ	73.7	0.47	0.47	73.7	73.7	0.73	73.7	73.7
	LP	66.3	0.36	0.32	61.3	88.5	0.72	44.2	88.5
	ET	74.5	0.50	0.49	80.0	65.5	0.72	83.6	65.5
PS2	RF	78.6	0.57	0.57	79.6	77.0	0.78	80.3	77.0
	ADB	73.7	0.47	0.47	73.7	75.4	0.74	72.1	75.4
	DT	59.8	0.20	0.19	65.0	42.6	0.51	77.0	42.6
	XGB	67.2	0.38	0.34	81.8	44.2	0.57	90.1	44.2
	QDA	50.8	0.02	0.01	53.8	11.4	0.18	90.1	11.4
	DJ	71.3	0.42	0.42	69.7	75.4	0.72	67.2	75.4
	LP	69.6	0.43	0.39	63.6	91.8	0.75	47.5	91.8
	ET (DHUpredET)	85.2	0.70	0.70	86.4	83.6	0.85	86.8	83.6

	RF	50.0	0.00	0.00	50.0	59.0	0.54	40.9	59.0
LSA	ADB	52.4	0.05	0.04	51.9	67.2	0.58	37.7	67.2
	DT	49.1	0.01	0.01	49.3	60.6	0.54	37.7	60.6
	XGB	50.0	0.00	0.00	0.00	0.00	0.00	100.0	0.00
	QDA	52.4	0.05	0.04	51.9	65.5	0.57	39.3	65.5
	DJ	56.5	0.14	0.13	54.7	75.4	0.63	37.7	75.4
	LP	50.0	0.00	0.00	50.0	68.8	0.57	31.1	68.8
	ET	47.5	0.05	0.04	48.0	60.6	0.53	34.4	60.6
	RF	67.6	0.34	0.34	67.8	65.5	0.66	68.8	65.5
BERT	ADB	58.2	0.16	0.16	57.3	63.9	0.60	52.4	63.9
	DT	54.1	0.08	0.08	53.8	57.3	0.55	50.8	57.3
	XGB	63.9	0.27	0.27	63.9	63.9	0.63	63.9	63.9
	QDA	50.8	0.09	0.01	1.00	1.64	0.03	100.0	0.01
	DJ	60.6	0.21	0.21	60.6	60.6	0.60	60.6	60.6
	LP	67.2	0.34	0.34	64.7	75.4	0.69	59.0	75.4
	ET	70.4	0.41	0.40	71.1	68.8	0.70	72.1	68.8
	RF	54.9	0.09	0.09	55.6	49.1	0.52	60.6	49.1
FastText	ADB	42.6	0.14	0.14	43.2	47.5	0.45	37.7	47.5
	DT	52.4	0.05	0.04	51.7	72.1	0.60	32.7	72.1
	XGB	50.0	0.00	0.00	0.00	0.00	0.00	100.0	0.00
	QDA	44.2	0.11	0.11	44.6	47.5	0.46	40.9	47.5
	DJ	53.2	0.07	0.06	52.2	75.4	0.61	31.1	75.4
	LP	50.0	0.00	0.00	50.0	40.9	0.45	59.0	40.9
	ET	46.7	0.06	0.06	46.7	47.5	0.47	45.9	47.5
	RF	70.6	0.41	0.41	70.9	71.1	0.71	70.4	70.1
Doc2Vec	ADB	63.1	0.26	0.26	61.7	68.8	0.65	57.3	68.8
	DT	60.6	0.21	0.21	61.0	59.0	0.60	62.3	59.0
	XGB	50.0	0.00	0.00	0.00	0.00	0.00	100.0	0.00
	QDA	50.0	0.00	0.00	0.00	0.00	0.00	100.0	0.00
	DJ	60.6	0.21	0.21	58.9	70.4	0.64	50.8	70.4
	LP	68.8	0.37	0.37	67.1	73.7	0.70	63.9	73.7
	ET	71.3	0.42	0.42	70.9	72.1	0.71	70.4	72.1
	RF	73.7	0.47	0.47	73.0	75.4	0.74	72.1	75.4
PseEIIP	ADB	65.5	0.31	0.31	64.1	70.4	0.67	60.6	70.4
	DT	59.8	0.19	0.19	60.3	57.3	0.58	62.3	57.3
	XGB	64.6	0.29	0.29	64.0	67.2	0.65	62.3	67.2
	QDA	63.9	0.27	0.27	63.0	67.2	0.65	60.6	67.2
	DJ	66.3	0.33	0.32	64.7	72.1	0.68	60.6	72.1
	LP	72.1	0.44	0.44	69.0	80.0	0.74	63.9	80.3
	ET	73.7	0.47	0.47	73.0	75.0	0.74	72.1	75.4

319

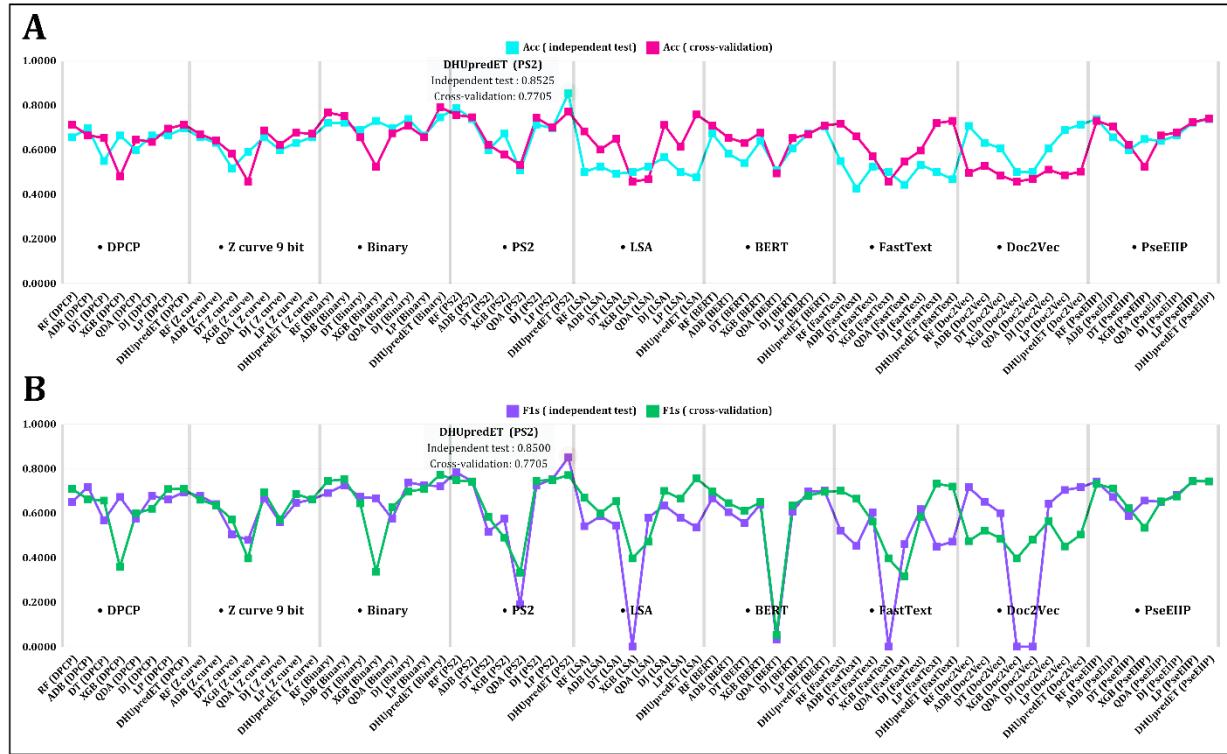
320

321 As seen from **Table 3**, most of the applied models have delivered promising results from various
 322 extracted features. Among them, ET using PS2 achieved the best results. Therefore, ET as a
 323 classifier and PS2 as the feature encoding method used to build DHUpredET model. DHUpredET
 324 obtains 85.3 rate of accuracy with excellent performances in other evaluation criteria. DHUpredET
 325 enhances the prediction accuracy by 10% compared to other physicochemical-based encoding
 326 methods. In NLP-based embedding, ADB and RF models acquired optimal results using LSA and
 327 FastText features. Overall, evaluations, the physicochemical-based PS2 features would be the most
 328 optimal features.

329

330

331 **Figure 4** shows the results of our comparison study using both the independent test method and
 332 the 5-fold cross-validation procedure. Our analysis is focused on comparing F1 scores and
 333 accuracy metrics. As shown in this figure, the results achieved for the independent test set and 5-
 334 fold cross-validation are consistent which demonstrates the generality and the robustness of our
 335 proposed method. DHUpredET consistently achieves excellent F1s and Acc across both
 336 assessment methodologies. This implies that DHUpredET has been effective at reducing both false
 337 positives and false negatives, indicating its effectiveness for its intended application. In the process
 338 of obtaining excellent scores in both measures across many evaluation techniques (independent
 339 test and cross-validation), DHUpredET demonstrates stability and durability in performance
 340 evaluation.



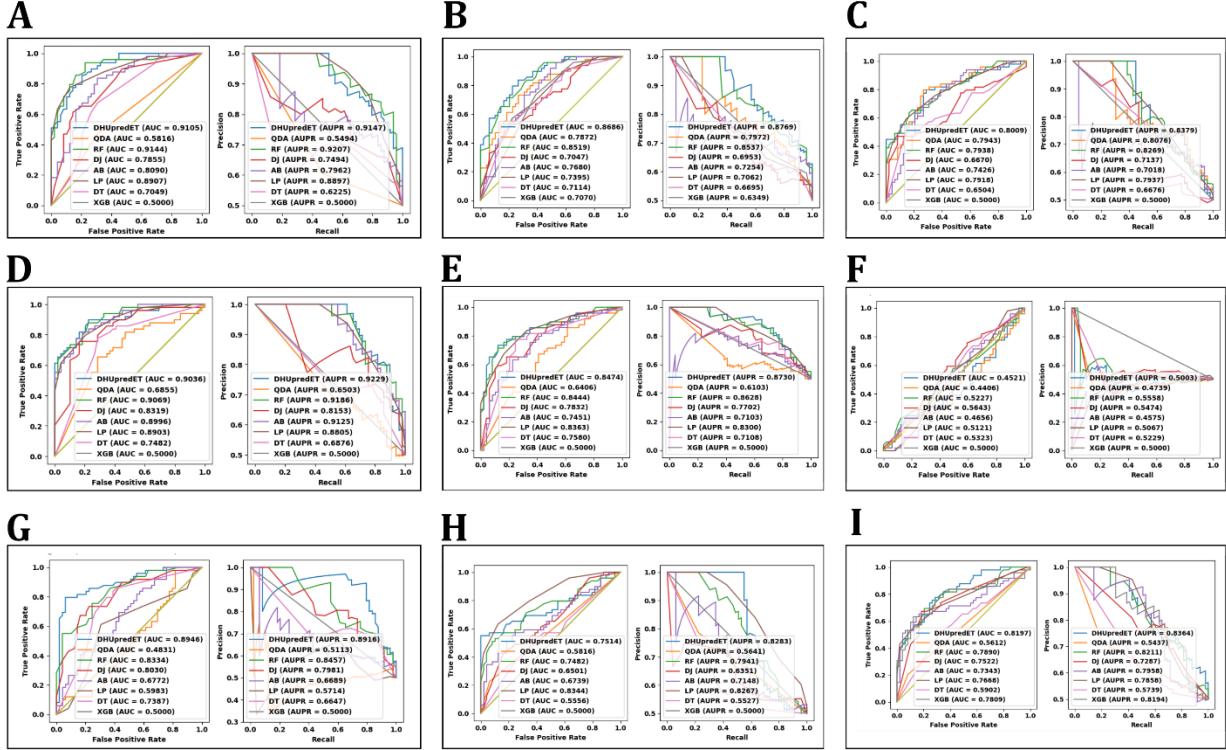
341
 342 **Figure 4.** Overall comparison of Accuracy and F1 score between independent test and 5-fold-
 343 cross-validation approach, where (A) Accuracy (Acc) (B) F1 Score (F1s)

344 The receiver operating characteristic (ROC) curve and precision-recall (PR) curve are two
 345 commonly used evaluation methods in the data science fields [61, 62]. The ROC curve illustrates
 346 the relationship between Sen and Spe, while the PR curve provides information about a model's
 347 efficiency, particularly when dealing with unbalanced datasets in which the number of negative
 348 cases significantly exceeds the positive ones. **Figure 5** shows the ROC and PR curves for our
 349 models using the independent test. Notably, the DHUpredET model performs adequately in all the
 350 subplots (A-I). In Subplot A (PS2, containing ROC and PR curves), DHUpredET receives
 351 noteworthy scores of 91.1% in ROC and 91.5% in PR curves. These results highlight that the
 352 model appropriately identifies both positive and negative categories, excelling at detecting positive
 353 examples while remaining specific for negative ones.
 354

355

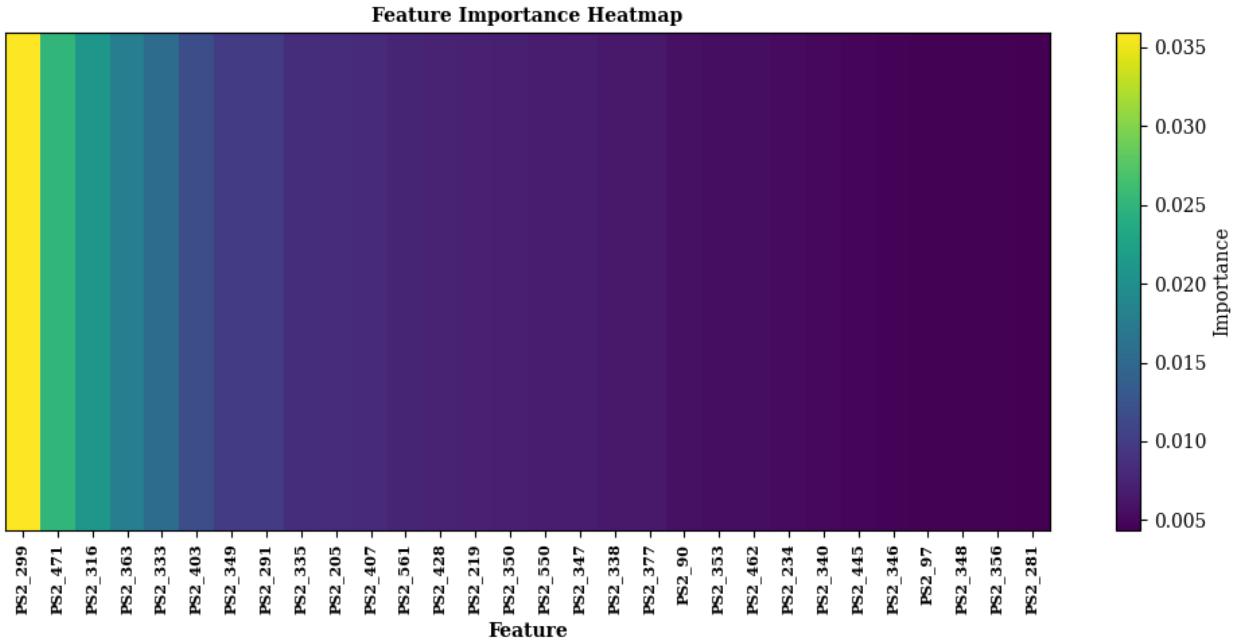
356

357 On the other hand, the PR curves demonstrates the model's outstanding performance. It also
 358 demonstrates the usefulness of the model in the presence of biased data sets. As shown in subplot
 359 D (Binary, containing ROC and PR curves), our model performed better with 92.3% of PR curves.
 360 Based on the findings of the present study, DHUpredET reliably performed the best across all
 361 subplots (A-I) compared to other models.



362
 363 **Figure 5.** ROC curves and PR curves visualization based on independent test method with all
 364 feature extraction methods. **(A)** PS2, containing ROC and PR curves **(B)** PseEIP, containing ROC
 365 and PR curves **(C)** DPCP, containing ROC and PR curves **(D)** Binary, containing ROC and PR
 366 curves **(E)** Z curve 9 bit, containing ROC and PR curves **(F)** LSA, containing ROC and PR curves
 367 **(G)** FastText, containing ROC and PR curves **(H)** Doc2Vec, containing ROC and PR curves **(I)**
 368 BERT, containing ROC and PR curves.
 369

370
 371 **3.1 Feature analysis**
 372 As our study found PS2 features have the most optimal attributes, therefore in **Figure 6**, We
 373 analyzed the contribution of each feature in our model's decision-making in PS2 feature
 374 extractions. Based on the results, we can conclude that the PS2_299 feature is the most important
 375 in our model, as indicated by its significant value achieved of 0.035. The PS2_471 feature follows
 376 closely after, with a notable significance score of 0.030, suggesting that it contributes considerably
 377 to the model's predictive capabilities. Furthermore, the PS2_316 feature appears as a significant
 378 contributor, instead with a significantly lower relevance score of 0.025. These feature significance
 379 ratings provide valuable details about our model's explanatory capacity. It is apparent that the three
 380 features—PS2_299, PS2_471, and PS2_316—play vital parts in the model's decision-making
 381 process. Their prominence emphasizes their importance in capturing the underlying patterns
 382 involved with D site identification.
 383



384
385 **Figure 6.** A heatmap of the DHUpredET model is based on PS2 features, and the study exhibits
386 the top 30 important features.
387

388 3.2 Compariosn of DHUpredET with the previous models

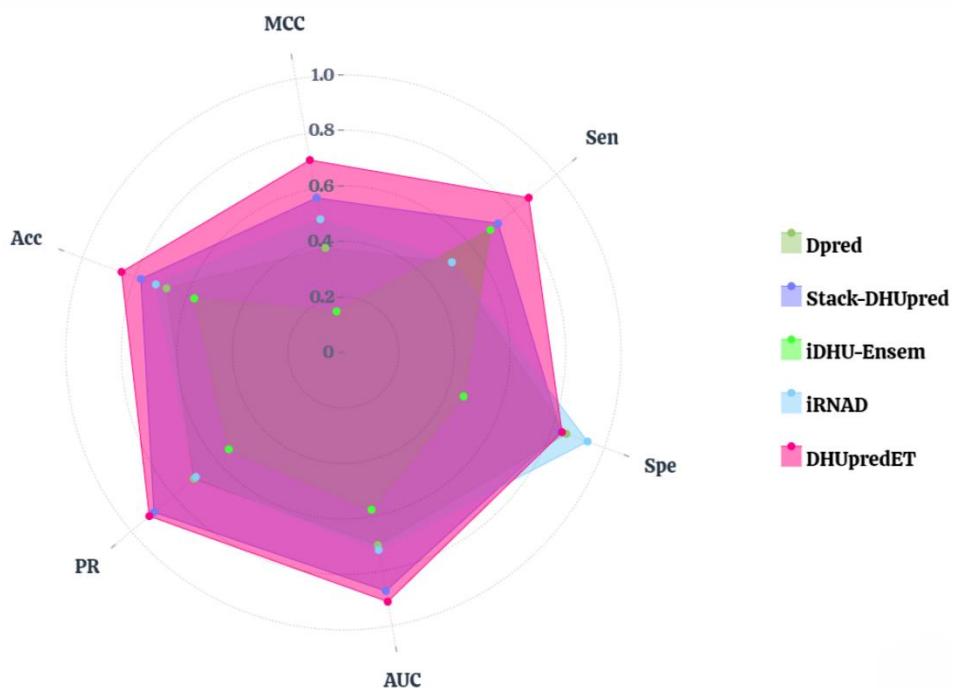
389
390 In **Table 4**, we compare DHUpredET to the state-of-the-art methods, evaluating their performance
391 across key metrics including Acc, Sen, AUC, MCC, PR, and Spe.
392

393 **Table 4.** Comparative analysis of the most recent methods based on the independent dataset with
394 DHUpredET model.
395

Model	Acc	MCC	Sen	Spe	AUC	PR
Dpred [9]	0.68	0.384	0.508	0.853	0.703	0.706
Stack-DHUpred [19]	0.778	0.567	0.725	0.830	0.871	0.893
iDHU-Ensem [18]	0.574	0.152	0.689	0.459	0.574	0.541
iRNAD [14]	0.721	0.489	0.508	0.934	0.721	0.696
DHUpredET	0.8525	0.7053	0.8689	0.8361	0.9105	0.9147

396 As shown in Table 4, DHUpredET achieves an accuracy of 85.3%, MCC of 70.53%, Sen of 56.9%,
397 AUC of 91.1%, and PR of 91.5%. These results greatly outperform previous methodologies.
398 DHUpredET model improved accuracy by more than 5% compared to all other models on the
399 independent test dataset. Despite the greatest increase in accuracy, our model performed
400 adequately across various evaluation metrics. The considerable improvement in MCC of more than
401 0.2 and Sen of more than 10% demonstrates that our model efficiently identifies positive classes
402 while maintaining a balanced Spe for negative classes. Additionally, the AUC score increased by
403 more than 3%, while the PR score boosted by 1%. The above results indicate that our model
404 outperforms state-of-the-art approaches and is well-suited to properly discriminating between
405 positive and negative classes in DHU prediction.
406

407 The comparison diagram in **Figure 8** visually illustrates the performance gap between the
 408 DHUpredET model and other methods. DHUpredET occupies a larger space, indicating its
 409 superior performance across multiple evaluation metrics compared to other methods..
 410



411
 412 **Figure 7.** Comparison of state-of-the-art models with DHUpredET based on accuracy, MCC,
 413 AUC, PR, specificity, and sensitivity.
 414
 415
 416

417 **4. Discussion**

418
 419 Dysregulation of RNA modifications has been associated with a variety of diseases, including
 420 malignancy and neurological disorders. Identifying DHU alteration sites may provide insight into
 421 disease causes and possible treatment targets. Experimental approaches for detecting RNA
 422 modifications, such as DHU, can be time-consuming and costly, with varying degrees of
 423 sensitivity and specificity.
 424

425 Computational techniques provide an alternative way for predicting modification sites with
 426 adequate accuracy, improving the prioritizing of experimental validation efforts. These approaches
 427 are important because they allow for the systematic analysis of massive datasets that would be
 428 hard to manage manually and provide a scalable answer to comprehending complicated biological
 429 processes by allowing for the full identification of D sites across a wide range of RNA molecules.
 430 The computer-based approach is especially useful when a preliminary high-throughput screening
 431 is required to focus future experimental efforts on the most promising targets. These predictions
 432 can thus accelerate the pace of discovery by directing empirical research towards the most relevant
 433 areas, enhancing the efficiency and effectiveness of biological study. In our extensive exploration
 434 of various studies, we observed a prevalent rhyme among authors, which presented different

435 techniques for detecting D sites. These representations generally employ various feature extraction
436 techniques and machine learning strategies. To conduct an exhaustive comparative analysis, we
437 utilized a variety of feature extraction approaches and analyzed 30 algorithms for each feature set.
438 After a comprehensive analysis, we evaluated the final models based on their performance and
439 selected the best features and methods. We explored integrating several feature sets, such as
440 physicochemical features only, compositional features only, and all features paired together, then
441 observed that PS2 features outperformed all of them. As a result, we adjusted our attention to using
442 specific factors, favoring both best-fit models and features. The aim was to develop a cost-effective
443 model without jeopardizing functionality. After rigorous testing and hyperparameter optimization,
444 we determined that the ET model was the best choice. Using PS2 features, which emphasize local
445 interdependence between nearby nucleotides, enables the model to improve its prediction
446 effectiveness by exploiting perplexing nucleotide interactions.

447

448

449

450 Conclusion

451 Dihydrouridine is a RNA modification that influences both the structural and functional patterns
452 of RNA, playing an essential role in manufacturing proteins and cellular adaption processes. Here
453 we developed a new machine learning technique, called DHUpredET, for predicting DHU using
454 extra tree as a classifier and PS2 feature encoding. Our proposed framework achieved 85.3% in
455 terms of prediction accuracy, significantly outperforming the state-of-the-art models found in the
456 literature by over 5%. Her ewe also comprehensively study variables that are strongly connected
457 to the identification of D sites. The DHUpredET model improves our ability to find dihydrouridine
458 sites with excellent accuracy and enhances our understanding of the underlying processes that
459 control RNA modifications from the computational perspective. The datasets used in the study
460 and the source code for DHUpredET are publicly available at <https://github.com/Shazzad-Shaon3404/DHUpredET-DHU-prediction.git>.

461

462

463

464 **Data and code availability**

465 The datasets of the study and the codes are available at this GitHub link:

466 <https://github.com/Shazzad-Shaon3404/DHUpredET-DHU-prediction.git>

467

468

469 **Contributions**

470 **Conceptualization:** Md. Fahim Sultan, Tasmin Karim, Md. Shazzad Hossain Shaon; **Data**
 471 **curation:** Tasmin Karim, Md. Shazzad Hossain Shaon, Md. Fahim Sultan, Md, Yeasin Biplob,
 472 Md. Shoaib Hossain Alshan; **Formal analysis:** Md. Fahim Sultan, Tasmin Karim, Md. Shazzad
 473 Hossain Shaon, Mst Shapna Akter; **Investigation:** Sayed Mehedi Azim, Md. Fahim Sultan,
 474 Tasmin Karim, Md. Shazzad Hossain Shaon, Mst Shapna Akter; **Methodology:** Md. Shazzad
 475 Hossain Shaon, Md. Fahim Sultan, Tasmin Karim, Sayed Mehedi Azim, Md. Shoaib Hossain
 476 Alshan, Mst Shapna Akter; **Project administration:** Md. Shazzad Hossain Shaon, Md. Fahim
 477 Sultan, Tasmin Karim, Sayed Mehedi Azim, Mst Shapna Akter, Iman Dehezangi; **Resources:** Md.
 478 Shazzad Hossain Shaon, Tasmin Karim, Sayed Mehedi Azim, Md. Yeasin Biplob, Md. Shoaib
 479 Hossain Alshan; **Supervision:** Iman Dehzangi ,Sayed Mehedi Azim , Md. Shazzad Hossain
 480 Shaon, Mst Shapna Akter; **Validation:** Iman Dehzangi ,Sayed Mehedi Azim , Md. Shazzad
 481 Hossain Shaon, Mst Shapna Akter; **Visualization:** Tasmin Karim, Md. Shazzad Hossain Shaon,
 482 Md. Shoaib Hossain Alshan; **Writing—original draft:** Md. Fahim Sultan, Tasmin Karim, Md.
 483 Shazzad Hossain Shaon, Md. Yeasin Biplob, Md. Shoaib Hossain Alshan; **Writing—review**
 484 **editing:** Iman Dehzangi, Sayed Mehedi Azim, Md. Shazzad Hossain Shaon, Mst Shapna Akter,
 485 Tasmin Karim; **Software:** Tasmin Karim, Md. Fahim Sultan, Md. Shazzad Hossain Shaon, Iman
 486 Dehezangi, Mst Shapna Akter;

487

488 The final version of the manuscript has been read and approved by all authors.

489

490 **References**

- 491 1. Post-translational modification [https://en.wikipedia.org/wiki/Post-](https://en.wikipedia.org/wiki/Post-translational_modification)
 492 [translational_modification](#). (accessed on 10 February 2024)
- 493 2. Arzumanian, V.A., Dolgalev, G.V., Kurbatov, I.Y., Kiseleva, O.I. and Poverennaya, E.V.,
 494 2022. Epitranscriptome: review of top 25 most-studied RNA modifications. International
 495 Journal of Molecular Sciences, 23(22), p.13851. <https://doi.org/10.3390/ijms232213851>
- 496 3. Boccaletto, P., Stefaniak, F., Ray, A., Cappannini, A., Mukherjee, S., Purta, E.,
 497 Kurkowska, M., Shirvanizadeh, N., Destefanis, E., Groza, P. and Avşar, G., 2022.
 498 MODOMICS: a database of RNA modification pathways. 2021 update. Nucleic acids
 499 research, 50(D1), pp.D231-D235. <https://doi.org/10.1093/nar/gkj084>
- 500 4. Cappannini, A., Ray, A., Purta, E., Mukherjee, S., Boccaletto, P., Moafinejad, S.N.,
 501 Lechner, A., Barchet, C., Klaholz, B.P., Stefaniak, F. and Bujnicki, J.M., 2024.
 502 MODOMICS: a database of RNA modifications and related information. 2023 update.
 503 Nucleic Acids Research, 52(D1), pp.D239-D244. <https://doi.org/10.1093/nar/gkad1083>
- 504 5. Shi, H., Wei, J. and He, C., 2019. Where, when, and how: context-dependent functions of
 505 RNA methylation writers, readers, and erasers. *Molecular cell*, 74(4), pp.640-650.
<https://doi.org/10.1016/j.molcel.2019.04.025>

- 508 6. Foltan, J.S., 2008. tRNA genes and the genetic code. *Journal of theoretical biology*, 253(3),
509 pp.469-482. <https://doi.org/10.1016/j.jtbi.2008.03.006>
- 510 7. Goodenbour, J.M. and Pan, T., 2006. Diversity of tRNA genes in eukaryotes. *Nucleic acids
511 research*, 34(21), pp.6137-6146. <https://doi.org/10.1093/nar/gkl725>
- 512 8. Edmonds, C.G., Crain, P.F., Gupta, R., Hashizume, T., Hocart, C.H., Kowalak, J.A.,
513 Pomerantz, S.C., Stetter, K.O. and McCloskey, J.A., 1991. Posttranscriptional
514 modification of tRNA in thermophilic archaea (Archaeabacteria). *Journal of bacteriology*,
515 173(10), pp.3138-3148. <https://doi.org/10.1128/jb.173.10.3138-3148.1991>
- 516 9. Yu, F., Tanaka, Y., Yamashita, K., Suzuki, T., Nakamura, A., Hirano, N., Suzuki, T., Yao,
517 M. and Tanaka, I., 2011. Molecular basis of dihydrouridine formation on tRNA.
518 *Proceedings of the National Academy of Sciences*, 108(49), pp.19593-19598.
519 <https://doi.org/10.1073/pnas.1112352108>
- 520 10. Draycott, A.S., Schaening-Burgos, C., Rojas-Duran, M.F., Wilson, L., Schärfen, L.,
521 Neugebauer, K.M., Nachtergael, S. and Gilbert, W.V., 2022. Transcriptome-wide
522 mapping reveals a diverse dihydrouridine landscape including mRNA. *PLoS Biology*,
523 20(5), p.e3001622. <https://doi.org/10.1371/journal.pbio.3001622>
- 524 11. Kato, T., Daigo, Y., Hayama, S., Ishikawa, N., Yamabuki, T., Ito, T., Miyamoto, M.,
525 Kondo, S. and Nakamura, Y., 2005. A novel human tRNA-dihydrouridine synthase
526 involved in pulmonary carcinogenesis. *Cancer research*, 65(13), pp.5638-5646.
527 <https://doi.org/10.1158/0008-5472.CAN-05-0600>
- 528 12. Mendez, M.F., 2017. Early-onset Alzheimer disease. *Neurologic clinics*, 35(2), pp.263-
529 281. <https://doi.org/10.1016/j.ncl.2017.01.005>
- 530 13. Durr, A., Gargiulo, M. and Feingold, J., 2012. The presymptomatic phase of Huntington
531 disease. *Revue neurologique*, 168(11), pp.806-808.
532 <https://doi.org/10.1016/j.neurol.2012.07.003>
- 533 14. Xu, Z.C., Feng, P.M., Yang, H., Qiu, W.R., Chen, W. and Lin, H., 2019. iRNAD: a
534 computational tool for identifying D modification sites in RNA sequence. *Bioinformatics*,
535 35(23), pp.4922-4929. <https://doi.org/10.1093/bioinformatics/btz358>
- 536 15. Dou, L., Zhou, W., Zhang, L., Xu, L. and Han, K., 2021. Accurate identification of RNA
537 D modification using multiple features. *RNA biology*, 18(12), pp.2236-2246.
538 <https://doi.org/10.1080/15476286.2021.1898160>
- 539 16. Zhu, H., Ao, C.Y., Ding, Y.J., Hao, H.X. and Yu, L., 2022. Identification of D Modification
540 sites using a random forest model based on nucleotide chemical properties. *International
541 Journal of Molecular Sciences*, 23(6), p.3044. <https://doi.org/10.3390/ijms23063044>
- 542 17. Suleman, M.T., Alkhailah, T., Alturise, F. and Khan, Y.D., 2022. DHU-Pred: accurate
543 prediction of dihydrouridine sites using position and composition variant features on
544 diverse classifiers. *PeerJ*, 10, p.e14104. <https://doi.org/10.7717/peerj.14104>
- 545 18. Suleman, M.T., Alturise, F., Alkhailah, T. and Khan, Y.D., 2023. iDHU-Ensem:
546 Identification of dihydrouridine sites through ensemble learning models. *Digital Health*, 9,
547 p.20552076231165963. <https://doi.org/10.1177/20552076231165963>
- 548 19. Harun-Or-Roshid, M., Maeda, K., Manavalan, B. and Kurata, H., 2024. Stack-DHUpred:
549 Advancing the accuracy of dihydrouridine modification sites detection via stacking
550 approach. *Computers in Biology and Medicine*, 169, p.107848.
551 <https://doi.org/10.1016/j.combiomed.2023.107848>
- 552 20. Sun, W.J., Li, J.H., Liu, S., Wu, J., Zhou, H., Qu, L.H. and Yang, J.H., 2016. RMBBase: a
553 resource for decoding the landscape of RNA modifications from high-throughput

- 554 sequencing data. Nucleic acids research, 44(D1), pp.D259-D265.
 555 <https://doi.org/10.1093/nar/gkv1036>
- 556 21. Fu, L., Niu, B., Zhu, Z., Wu, S. and Li, W., 2012. CD-HIT: accelerated for clustering the
 557 next-generation sequencing data. Bioinformatics, 28(23), pp.3150-3152.
 558 <https://doi.org/10.1093/bioinformatics/bts565>
- 559 22. Manavalan, B., Hasan, M.M., Basith, S., Gosu, V., Shin, T.H. and Lee, G., 2020. Empirical
 560 comparison and analysis of web-based DNA N4-methylcytosine site prediction tools.
 561 Molecular Therapy-Nucleic Acids, 22, pp.406-420.
 562 <https://doi.org/10.1016/j.omtn.2020.09.010>
- 563 23. Hasan, M.M., Tsukiyama, S., Cho, J.Y., Kurata, H., Alam, M.A., Liu, X., Manavalan, B.
 564 and Deng, H.W., 2022. Deepm5C: a deep-learning-based hybrid framework for identifying
 565 human RNA N5-methylcytosine sites using a stacking strategy. Molecular Therapy, 30(8),
 566 pp.2856-2867. <https://doi.org/10.1016/j.ymthe.2022.05.001>
- 567 24. Hasan, M.M., Basith, S., Khatun, M.S., Lee, G., Manavalan, B. and Kurata, H., 2021.
 568 Meta-i6mA: an interspecies predictor for identifying DNA N 6-methyladenine sites of
 569 plant genomes by exploiting informative features in an integrative machine-learning
 570 framework. Briefings in Bioinformatics, 22(3), p.bbbaa202.
 571 <https://doi.org/10.1093/bib/bbbaa202>
- 572 25. Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer,
 573 L. and Stoyanov, V., 2019. Roberta: A robustly optimized bert pretraining approach. arXiv
 574 preprint arXiv:1907.11692. <https://doi.org/10.48550/arXiv.1907.11692>
- 575 26. Devlin, J., Chang, M.W., Lee, K. and Toutanova, K., 2018. Bert: Pre-training of deep
 576 bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.
 577 <https://doi.org/10.48550/arXiv.1810.04805>
- 578 27. Huo, H. and Iwaihara, M., 2020. Utilizing BERT pretrained models with various fine-tune
 579 methods for subjectivity detection. In Web and Big Data: 4th International Joint
 580 Conference, APWeb-WAIM 2020, Tianjin, China, September 18-20, 2020, Proceedings,
 581 Part II 4 (pp. 270-284). Springer International Publishing. https://doi.org/10.1007/978-3-030-60290-1_21
- 583 28. Su, W., Zhu, X., Cao, Y., Li, B., Lu, L., Wei, F. and Dai, J., 2019. ViLbert: Pre-training of
 584 generic visual-linguistic representations. arXiv preprint arXiv:1908.08530.
 585 <https://doi.org/10.48550/arXiv.1908.08530>
- 586 29. Bojanowski, P., Grave, E., Joulin, A. and Mikolov, T., 2017. Enriching word vectors with
 587 subword information. Transactions of the association for computational linguistics, 5,
 588 pp.135-146. https://doi.org/10.1162/tacl_a_00051
- 589 30. Choi, J. and Lee, S.W., 2020. Improving FastText with inverse document frequency of
 590 subwords. Pattern Recognition Letters, 133, pp.165-172.
 591 <https://doi.org/10.1016/j.patrec.2020.03.003>
- 592 31. Kwon, H., Kim, J. and Park, Y., 2017. Applying LSA text mining technique in envisioning
 593 social impacts of emerging technologies: The case of drone technology. Technovation, 60,
 594 pp.15-28. <https://doi.org/10.1016/j.technovation.2017.01.001>
- 595 32. Yu, B., Xu, Z.B. and Li, C.H., 2008. Latent semantic analysis for text categorization using
 596 neural network. Knowledge-Based Systems, 21(8), pp.900-904.
 597 <https://doi.org/10.1016/j.knosys.2008.03.045>

- 598 33. Chen, W.K., Chen, L.S. and Pan, Y.T., 2021. A text mining-based framework to discover
599 the important factors in text reviews for predicting the views of live streaming. Applied
600 Soft Computing, 111, p.107704. <https://doi.org/10.1016/j.asoc.2021.107704>
- 601 34. Le, Q. and Mikolov, T., 2014, June. Distributed representations of sentences and
602 documents. In International conference on machine learning (pp. 1188-1196). PMLR.
- 603 35. Chen, Q. and Sokolova, M., 2021. Specialists, scientists, and sentiments: Word2Vec and
604 Doc2Vec in analysis of scientific and medical texts. SN Computer Science, 2, pp.1-11.
605 <https://doi.org/10.1007/s42979-021-00807-1>
- 606 36. Mishra, S., Pappu, A. and Bhamidipati, N., 2019, May. Inferring advertiser sentiment in
607 online articles using wikipedia footnotes. In Companion proceedings of the 2019 world
608 wide web conference (pp. 1224-1231). <https://doi.org/10.1145/3308560.3316752>
- 609 37. nidia, Robson P., et al. "Feature extraction approaches for biological sequences: a
610 comparative study of mathematical features." Briefings in Bioinformatics 22.5 (2021):
611 bbab011. <https://doi.org/10.1093/bib/bbab011>
- 612 38. Han, Siyu, et al. "LncFinder: an integrated platform for long non-coding RNA
613 identification utilizing sequence intrinsic composition, structural information and
614 physicochemical property." Briefings in bioinformatics 20.6 (2019): 2009-2027.
615 <https://doi.org/10.1093/bib/bby065>
- 616 39. Manavalan, B., Basith, S., Shin, T.H., Lee, D.Y., Wei, L. and Lee, G., 2019. 4mCpred-EL:
617 an ensemble learning framework for identification of DNA N4-methylcytosine sites in the
618 mouse genome. Cells, 8(11), p.1332. <https://doi.org/10.3390/cells8111332>
- 619 40. Chen, Z., Chen, Y.Z., Wang, X.F., Wang, C., Yan, R.X. and Zhang, Z., 2011. Prediction
620 of ubiquitination sites by using the composition of k-spaced amino acid pairs. PloS one,
621 6(7), p.e22930. <https://doi.org/10.1371/journal.pone.0022930>
- 622 41. Chen, Z., Zhou, Y., Song, J. and Zhang, Z., 2013. hCKSAAP_UbSite: improved prediction
623 of human ubiquitination sites by exploiting amino acid pattern and properties. Biochimica
624 et Biophysica Acta (BBA)-Proteins and Proteomics, 1834(8), pp.1461-1467.
625 <https://doi.org/10.1016/j.bbapap.2013.04.006>
- 626 42. Liu, B., Gao, X. and Zhang, H., 2019. BioSeq-Analysis2. 0: an updated platform for
627 analyzing DNA, RNA and protein sequences at sequence level and residue level based on
628 machine learning approaches. Nucleic acids research, 47(20), pp.e127-e127.
629 <https://doi.org/10.1093/nar/gkz740>
- 630 43. Doench, J.G., Fusi, N., Sullender, M., Hegde, M., Vaimberg, E.W., Donovan, K.F., Smith,
631 I., Tothova, Z., Wilen, C., Orchard, R. and Virgin, H.W., 2016. Optimized sgRNA design
632 to maximize activity and minimize off-target effects of CRISPR-Cas9. Nature
633 biotechnology, 34(2), pp.184-191. <https://www.nature.com/articles/nbt.3437#citeas>
- 634 44. Gao, F. and Zhang, C.T., 2004. Comparison of various algorithms for recognizing short
635 coding sequences of human genes. Bioinformatics, 20(5), pp.673-681.
636 <https://doi.org/10.1093/bioinformatics/btg467>
- 637 45. Ying, C., Qi-Guang, M., Jia-Chen, L. and Lin, G., 2013. Advance and prospects of
638 AdaBoost algorithm. Acta Automatica Sinica, 39(6), pp.745-758.
639 [https://doi.org/10.1016/S1874-1029\(13\)60052-X](https://doi.org/10.1016/S1874-1029(13)60052-X)
- 640 46. Garza, S.E. and Schaeffer, S.E., 2019. Community detection with the label propagation
641 algorithm: a survey. Physica A: Statistical Mechanics and its Applications, 534, p.122058.
642 <https://doi.org/10.1016/j.physa.2019.122058>

- 643 47. Tharwat, A., 2016. Linear vs. quadratic discriminant analysis classifier: a tutorial.
644 International Journal of Applied Pattern Recognition, 3(2), pp.145-180.
645 <https://doi.org/10.1504/IJAPR.2016.079050>
- 646 48. Bose, S., Pal, A., SahaRay, R. and Nayak, J., 2015. Generalized quadratic discriminant
647 analysis. Pattern Recognition, 48(8), pp.2676-2684.
648 <https://doi.org/10.1016/j.patcog.2015.02.016>
- 649 49. Chen, T., He, T., Benesty, M., Khotilovich, V., Tang, Y., Cho, H., Chen, K., Mitchell, R.,
650 Cano, I. and Zhou, T., 2015. Xgboost: extreme gradient boosting. R package version 0.4-
651 2, 1(4), pp.1-4.
- 652 50. Chen, T. and Guestrin, C., 2016, August. Xgboost: A scalable tree boosting system. In
653 Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and
654 data mining (pp. 785-794). <https://doi.org/10.1145/2939672.2939785>
- 655 51. Suthaharan, S. and Suthaharan, S., 2016. Decision tree learning. Machine Learning Models
656 and Algorithms for Big Data Classification: Thinking with Examples for Effective
657 Learning, pp.237-269. https://doi.org/10.1007/978-1-4899-7641-3_10
- 658 52. Rigatti, S.J., 2017. Random forest. Journal of Insurance Medicine, 47(1), pp.31-39.
659 <https://doi.org/10.17849/insm-47-01-31-39.1>
- 660 53. Qi, Y., 2012. Random forest for bioinformatics. Ensemble machine learning: Methods and
661 applications, pp.307-323. https://doi.org/10.1007/978-1-4419-9326-7_11
- 662 54. Shotton, J., Sharp, T., Kohli, P., Nowozin, S., Winn, J. and Criminisi, A., 2013. Decision
663 jungles: Compact and rich models for classification. Advances in neural information
664 processing systems, 26.
- 665 55. Sharaff, A. and Gupta, H., 2019. Extra-tree classifier with metaheuristics approach for
666 email classification. In Advances in Computer Communication and Computational
667 Sciences: Proceedings of IC4S 2018 (pp. 189-197). Springer Singapore.
668 https://doi.org/10.1007/978-981-13-6861-5_17
- 669 56. Oostwal, E., Straat, M. and Biehl, M., 2021. Hidden unit specialization in layered neural
670 networks: ReLU vs. sigmoidal activation. Physica A: Statistical Mechanics and its
671 Applications, 564, p.125517. <https://doi.org/10.1016/j.physa.2020.125517>
- 672 57. Umakantha, N., 2016. A new approach to probability theory with reference to statistics and
673 statistical physics. Journal of Modern Physics, 7(09), p.989.
674 <http://dx.doi.org/10.4236/jmp.2016.79090>
- 675 58. Kraemer HC. Kappa coefficient. Wiley StatsRef: statistics reference online. 2014 Apr 14:1-
676 4. <https://doi.org/10.1002/9781118445112.stat00365.pub2>
- 677 59. Satu, M.S., Akter, T. and Uddin, M.J., 2017, February. Performance analysis of classifying
678 localization sites of protein using data mining techniques and artificial neural networks. In
679 2017 International Conference on Electrical, Computer and Communication Engineering
680 (ECCE) (pp. 860-865). IEEE. <https://doi.org/10.1109/ECACE.2017.7913023>
- 681 60. Thabtah, F., 2019. Machine learning in autistic spectrum disorder behavioral research: A
682 review and ways forward. Informatics for Health and Social Care, 44(3), pp.278-297.
683 <https://doi.org/10.1080/17538157.2017.1399132>
- 684 61. Hoo, Z.H., Candlish, J. and Teare, D., 2017. What is an ROC curve?. *Emergency Medicine
685 Journal*, 34(6), pp.357-359. <https://doi.org/10.1136/emermed-2017-206735>

- 686 62. Davis, J. and Goadrich, M., 2006, June. The relationship between Precision-Recall and
687 ROC curves. In *Proceedings of the 23rd international conference on Machine learning* (pp.
688 233-240). <https://doi.org/10.1145/1143844.1143874>
- 689

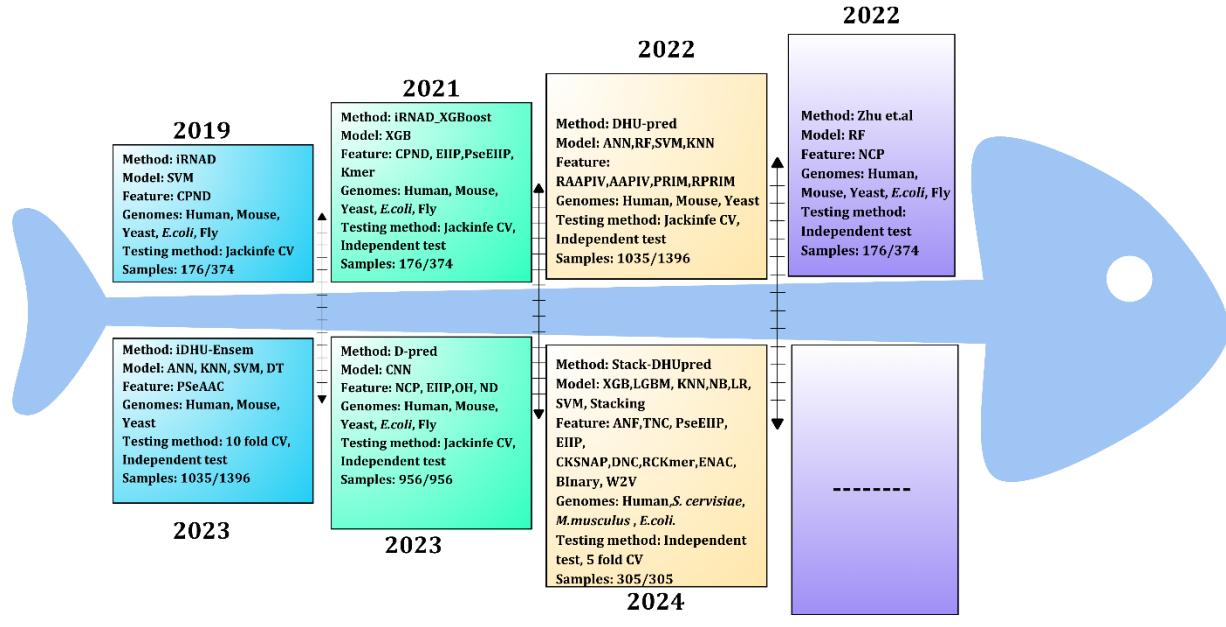


Figure 1. An overview of the existing methods in various years, represented with a fishbone diagram.

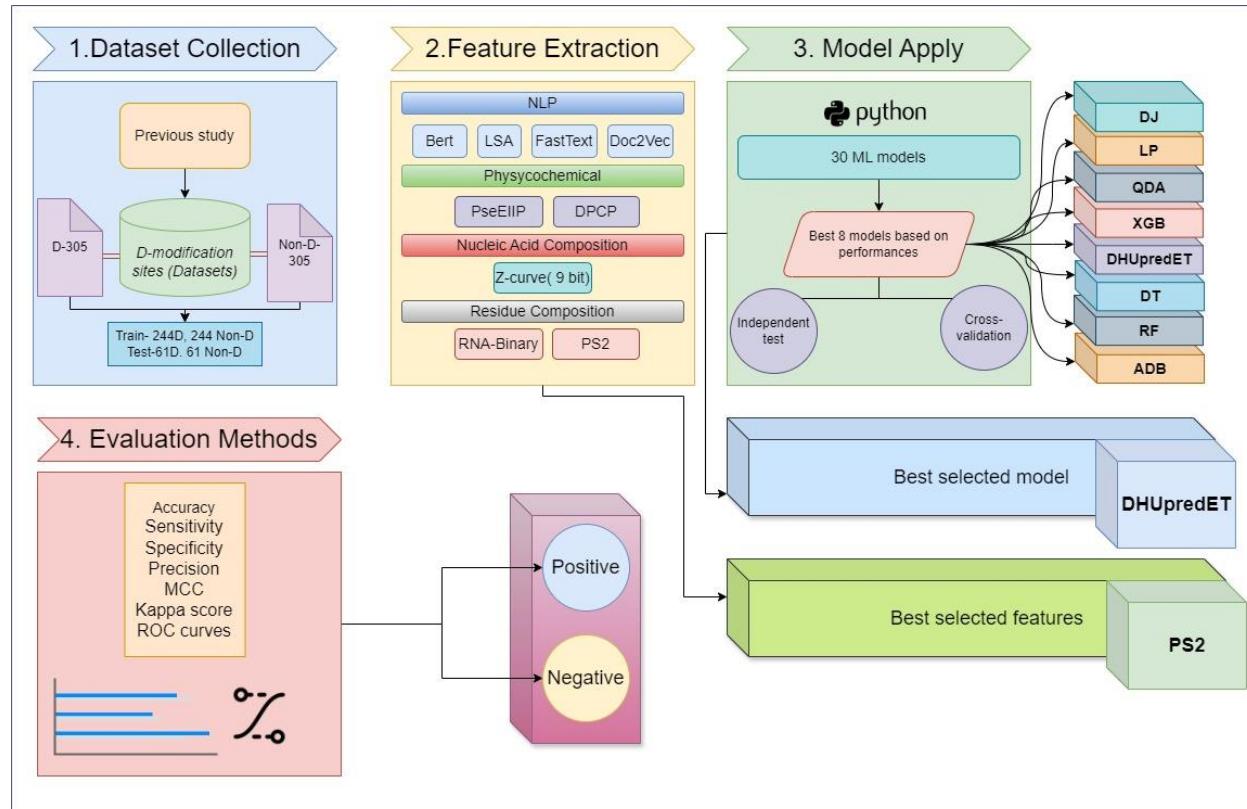


Figure 2. Workflows of the current study. Dataset collection from the previous studies, feature extraction approach of four kinds of descriptors with nine feature encoding methods, application of the various models, and selection process, the overall comparison of the models, and selecting the optimal feature and best-fit models for the study.

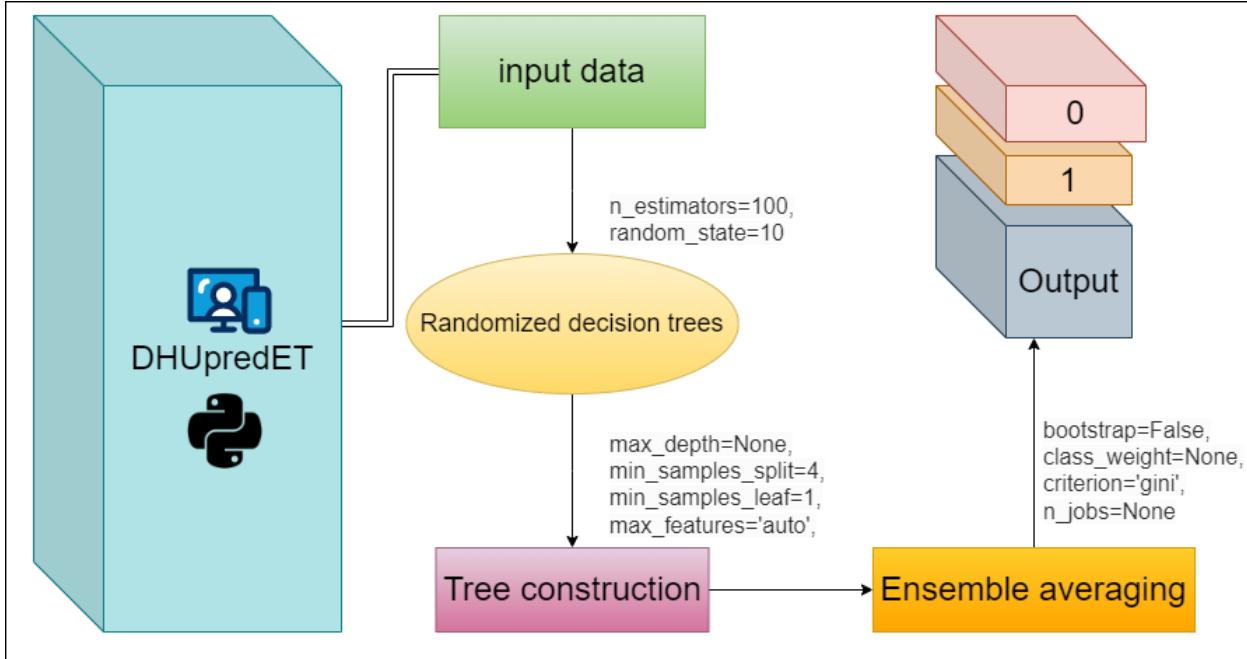


Figure 3. Overall strategies of the DHUpredET model.

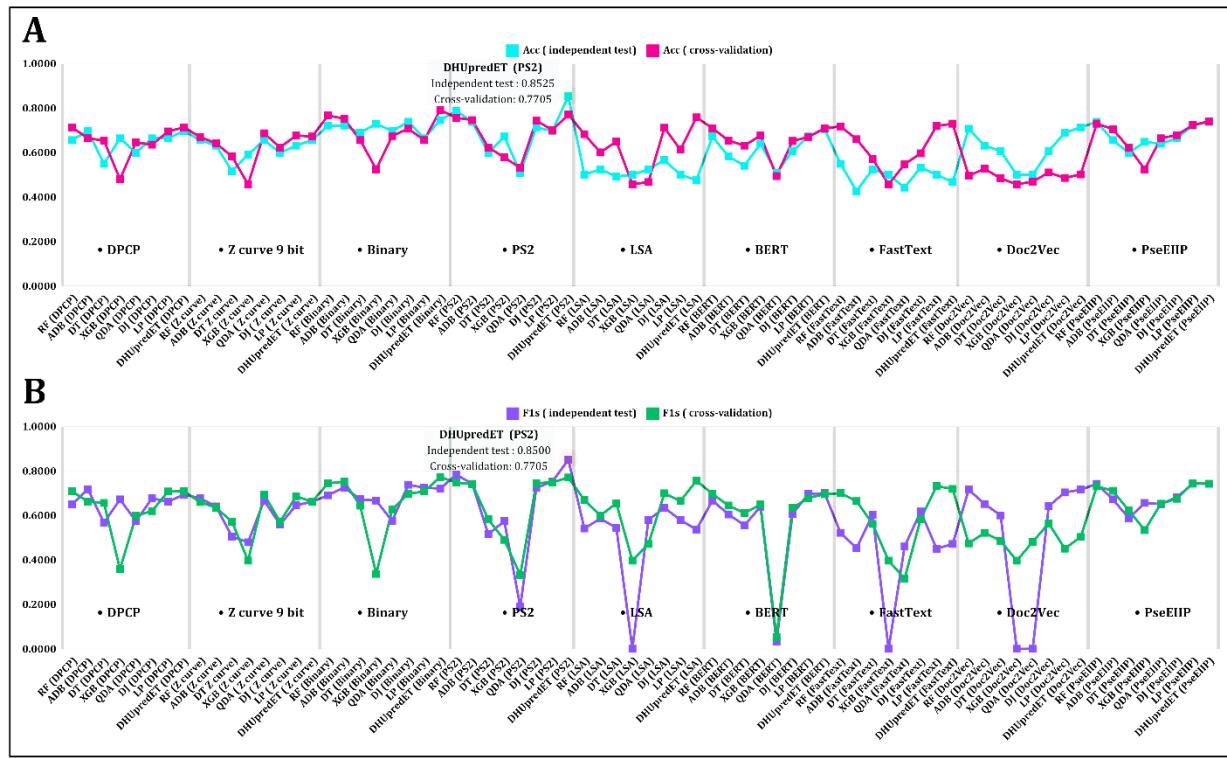


Figure 4. Overall comparison of Accuracy and F1 score between independent test and 5-fold-cross-validation approach, where **(A)** Accuracy (Acc) **(B)** F1 Score (F1s)

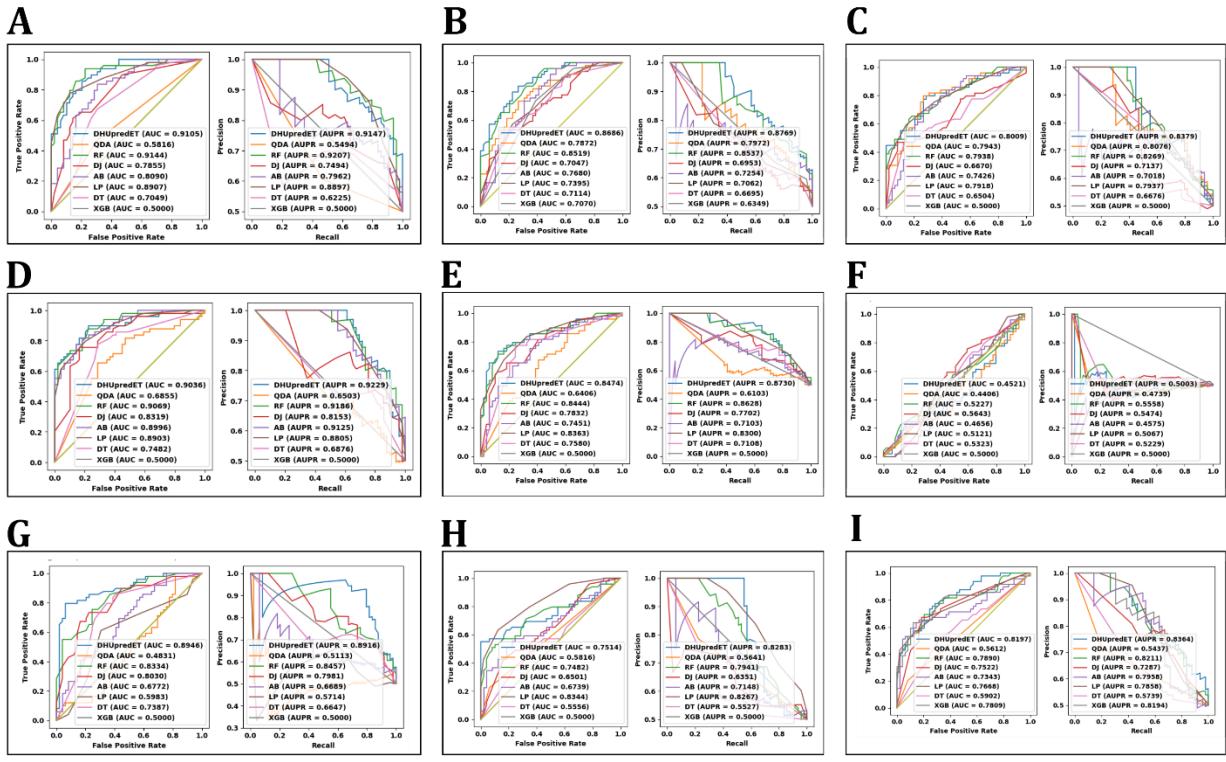


Figure 5. ROC curves and PR curves visualization based on independent test method with all feature extraction methods. **(A)** PS2, containing ROC and PR curves **(B)** PseEIP, containing ROC and PR curves **(C)** DPCP, containing ROC and PR curves **(D)** Binary, containing ROC and PR curves **(E)** Z curve 9 bit, containing ROC and PR curves **(F)** LSA, containing ROC and PR curves **(G)** FastText, containing ROC and PR curves **(H)** Doc2Vec, containing ROC and PR curves **(I)** BERT, containing ROC and PR curves.

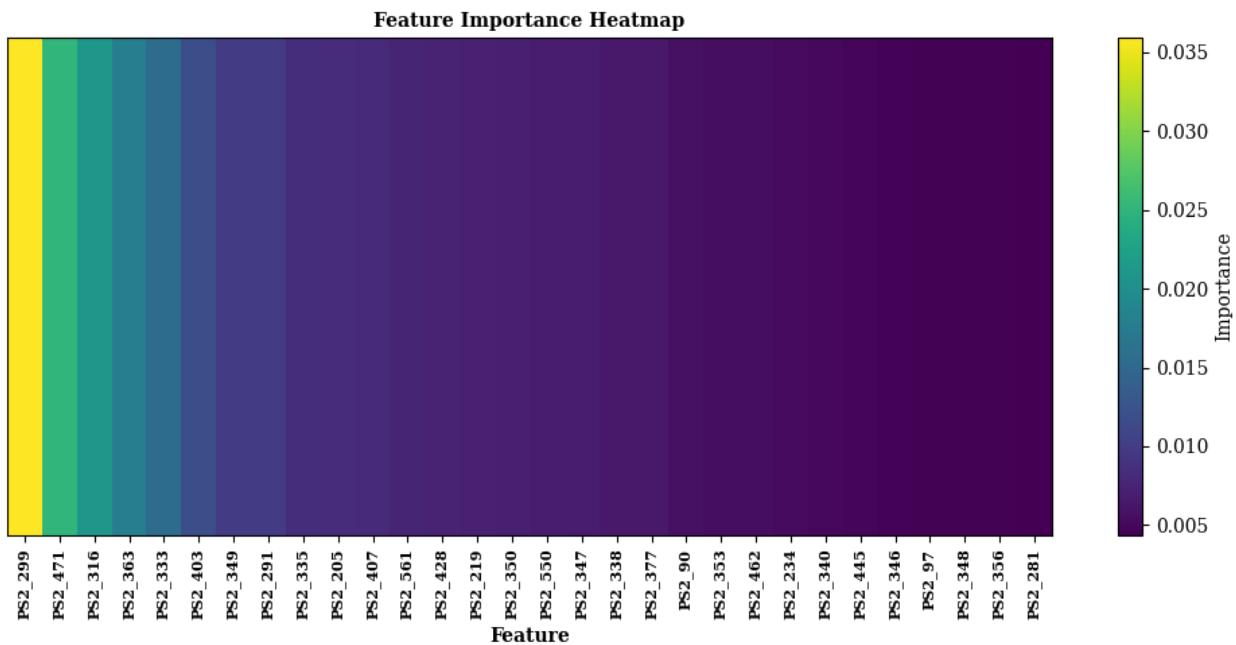


Figure 6. A heatmap of the DHUpredET model is based on PS2 features, and the study exhibits the top 30 important features.

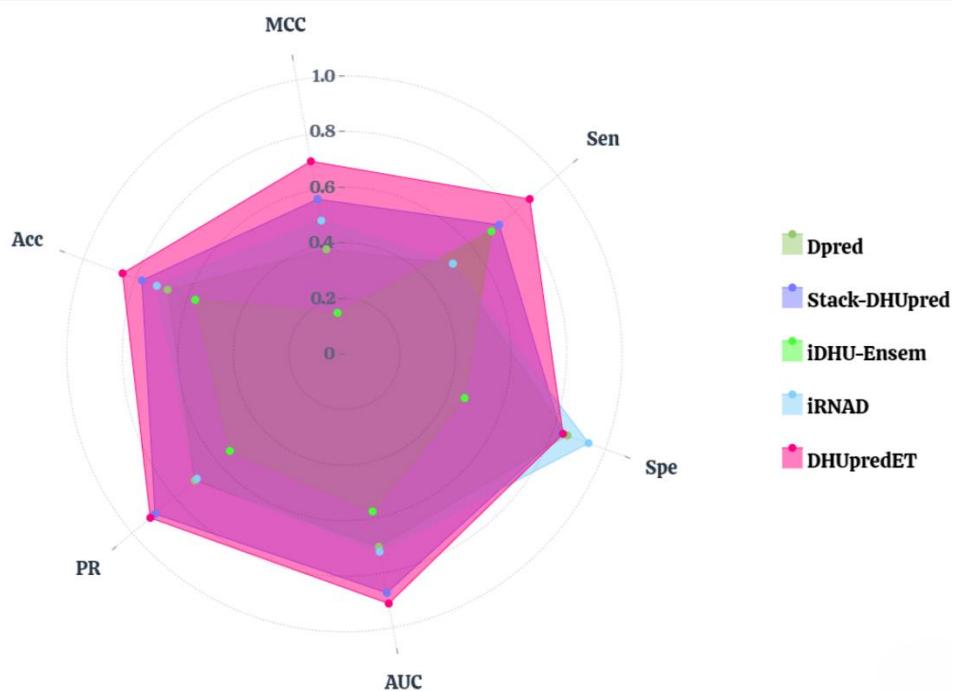


Figure 7. Comparison of state-of-the-art models with DHUpredET based on accuracy, MCC, AUC, PR, specificity, and sensitivity.



Click here to access/download
Supplementary Material
suppliments files (1).xlsx

