

FEASIBILITY TEST UPDATES

Resampling/SMOTE

Resampling and SMOTE are possible through the Imbalance Learning (*imblearn*) library (already added on CH3: Libraries). For Oversampling, imblearn's [RandomOverSampler](#) was used. For SMOTE, [SMOTEN](#) (Synthetic Minority Over-sampling Technique for Nominal) was used (SMOTE but designed to handle categorical data). Do note that resampling (regardless of the technique) may incur overfitting for Benign samples for the case the Oliveira dataset. Hence, there must be proof that can debunk or at least alleviate overfitting concerns which can be done through k-folds testing or model robustness test.

Recent [tests](#) on Oliveira (via HGBT Model) suggests that Oversampling is more 'equal' than SMOTE is due to the latter being skewed on certain benign samples. However, in terms of model performance results, the training dataset that had undergone SMOTE outperforms Oversampling.

Oversampling	Top 5 Most Repeated Samples																														
	<table><tr><th>Sample</th><th>Quantity in Resampled</th><th>% in Resampled</th></tr><tr><td>03384ab6368b68ed16ecb9e6352539af</td><td>90</td><td>0.22%</td></tr><tr><td>0822ec2ba98d291e5bfc836bc3686096</td><td>90</td><td>0.22%</td></tr><tr><td>f78ea80cec007b2c32fb10f9c6c82f39</td><td>88</td><td>0.21%</td></tr><tr><td>075323e77815ee8bcc7854ce23955a15</td><td>79</td><td>0.19%</td></tr><tr><td>79b78bb3d583748040c41ded09555fd3</td><td>72</td><td>0.17%</td></tr></table>	Sample	Quantity in Resampled	% in Resampled	03384ab6368b68ed16ecb9e6352539af	90	0.22%	0822ec2ba98d291e5bfc836bc3686096	90	0.22%	f78ea80cec007b2c32fb10f9c6c82f39	88	0.21%	075323e77815ee8bcc7854ce23955a15	79	0.19%	79b78bb3d583748040c41ded09555fd3	72	0.17%												
Sample	Quantity in Resampled	% in Resampled																													
03384ab6368b68ed16ecb9e6352539af	90	0.22%																													
0822ec2ba98d291e5bfc836bc3686096	90	0.22%																													
f78ea80cec007b2c32fb10f9c6c82f39	88	0.21%																													
075323e77815ee8bcc7854ce23955a15	79	0.19%																													
79b78bb3d583748040c41ded09555fd3	72	0.17%																													
	Model Performance with Oversampled Dataset																														
	<table><tr><th></th><th>precision</th><th>recall</th><th>f1-score</th><th>support</th></tr><tr><td>0</td><td>0.9110</td><td>0.8571</td><td>0.8832</td><td>12827</td></tr><tr><td>1</td><td>0.8653</td><td>0.9164</td><td>0.8901</td><td>12852</td></tr><tr><td>accuracy</td><td></td><td></td><td>0.8868</td><td>25679</td></tr><tr><td>macro avg</td><td>0.8882</td><td>0.8868</td><td>0.8867</td><td>25679</td></tr><tr><td>weighted avg</td><td>0.8881</td><td>0.8868</td><td>0.8867</td><td>25679</td></tr></table>		precision	recall	f1-score	support	0	0.9110	0.8571	0.8832	12827	1	0.8653	0.9164	0.8901	12852	accuracy			0.8868	25679	macro avg	0.8882	0.8868	0.8867	25679	weighted avg	0.8881	0.8868	0.8867	25679
	precision	recall	f1-score	support																											
0	0.9110	0.8571	0.8832	12827																											
1	0.8653	0.9164	0.8901	12852																											
accuracy			0.8868	25679																											
macro avg	0.8882	0.8868	0.8867	25679																											
weighted avg	0.8881	0.8868	0.8867	25679																											
SMOTEN (k_neighbors=5)	Top 5 Most Repeated Samples																														
	<table><tr><th>Sample</th><th>Quantity in Resampled</th><th>% in Resampled</th></tr><tr><td>0da6a786018d3e267de65f253277a1e0</td><td>5965</td><td>14.30%</td></tr><tr><td>302586218f78bb35439df31b54685ad0</td><td>3728</td><td>8.94%</td></tr><tr><td>0327301655f2e1c0b0bd4636a3349216</td><td>1895</td><td>4.54%</td></tr><tr><td>1707b149c4d1000242321167eec80eed</td><td>554</td><td>1.33%</td></tr><tr><td>3506356c329758e4f703cd2103d7daab</td><td>543</td><td>1.30%</td></tr></table>	Sample	Quantity in Resampled	% in Resampled	0da6a786018d3e267de65f253277a1e0	5965	14.30%	302586218f78bb35439df31b54685ad0	3728	8.94%	0327301655f2e1c0b0bd4636a3349216	1895	4.54%	1707b149c4d1000242321167eec80eed	554	1.33%	3506356c329758e4f703cd2103d7daab	543	1.30%												
Sample	Quantity in Resampled	% in Resampled																													
0da6a786018d3e267de65f253277a1e0	5965	14.30%																													
302586218f78bb35439df31b54685ad0	3728	8.94%																													
0327301655f2e1c0b0bd4636a3349216	1895	4.54%																													
1707b149c4d1000242321167eec80eed	554	1.33%																													
3506356c329758e4f703cd2103d7daab	543	1.30%																													
	Model Performance with SMOTEN Dataset																														
	<table><tr><th></th><th>precision</th><th>recall</th><th>f1-score</th><th>support</th></tr><tr><td>0</td><td>0.9376</td><td>0.9411</td><td>0.9393</td><td>12827</td></tr><tr><td>1</td><td>0.9410</td><td>0.9374</td><td>0.9392</td><td>12852</td></tr><tr><td>accuracy</td><td></td><td></td><td>0.9393</td><td>25679</td></tr><tr><td>macro avg</td><td>0.9393</td><td>0.9393</td><td>0.9393</td><td>25679</td></tr><tr><td>weighted avg</td><td>0.9393</td><td>0.9393</td><td>0.9393</td><td>25679</td></tr></table>		precision	recall	f1-score	support	0	0.9376	0.9411	0.9393	12827	1	0.9410	0.9374	0.9392	12852	accuracy			0.9393	25679	macro avg	0.9393	0.9393	0.9393	25679	weighted avg	0.9393	0.9393	0.9393	25679
	precision	recall	f1-score	support																											
0	0.9376	0.9411	0.9393	12827																											
1	0.9410	0.9374	0.9392	12852																											
accuracy			0.9393	25679																											
macro avg	0.9393	0.9393	0.9393	25679																											
weighted avg	0.9393	0.9393	0.9393	25679																											

Concerns on Model Robustness Test

While it is certain that it will be removed at its current definition (i.e., using different datasets), it is also possible that it will be retained, albeit being redefined in certain aspects. It is possible that part of Oliveira can serve as a stand-in for the external dataset.

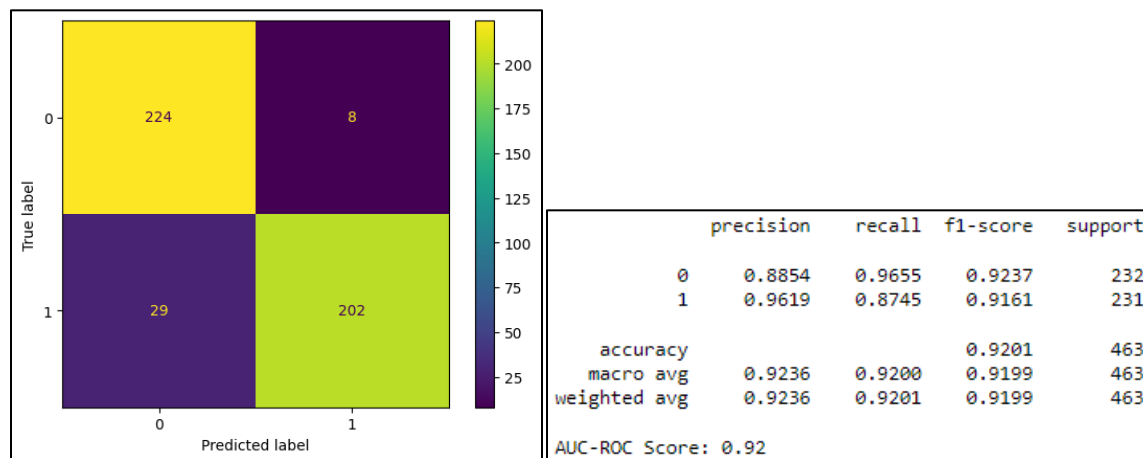
Recent [tests](#) suggests that it is indeed possible. But do note that Oliveira, at its core, is skewed for Malicious samples, hence it will never score $\sim 90\%$ for on certain metrics (i.e., precision and recall) on individual labels. The two images below show the 90:10 ratio for train (left) and reserve (right) splits respectively.

train_split shape: (39488, 102)	reserve_split shape: (4388, 102)
train_split value_counts:	malware
malware	1 4284
1 38513	0 104
0 975	Name: count, dtype: int64
Name: count, dtype: int64	

LightGBM

LightGBM was tested on MalbehavD-V1 to determine implementation on similar datasets (i.e., one with encoded/numerical data). The dataset (both train and input/test) must be encoded (e.g., LabelEncoded) before use as its implementation of 'support for categorical data' quite misleading as per its documentation. Recent [simple tests](#) also suggest that GPUs (tested using AMD GPU on Win11) are supported.

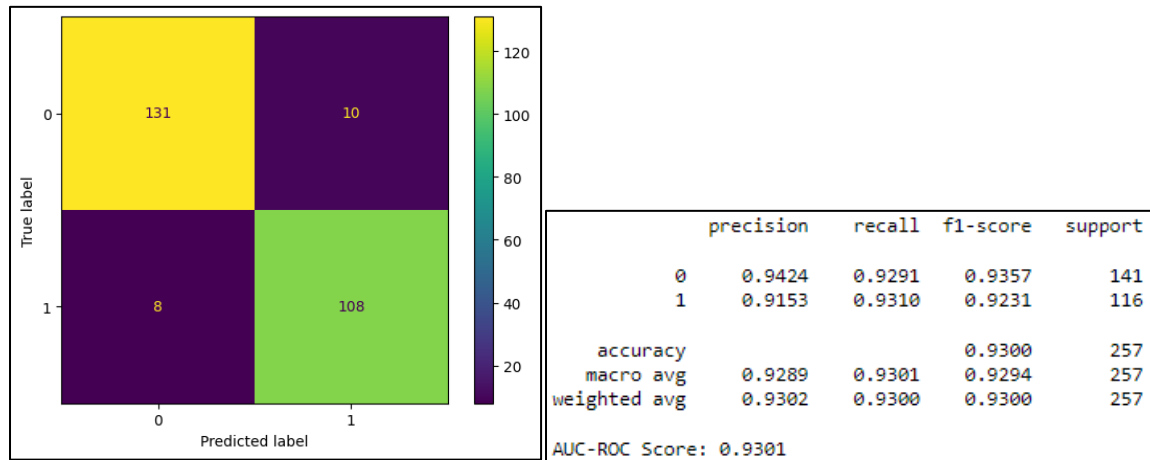
Split Sample Test Results:



K-Folds Test Results:

```
[['accuracy', 'f1_score', 'precision', 'recall', 'roc_auc', 'time'],
 [0.9071, 0.907, 0.9401, 0.8718, 0.9075, 25.8268],
 [0.9395, 0.9395, 0.9682, 0.9103, 0.9398, 23.8312],
 [0.9266, 0.9265, 0.963, 0.8889, 0.927, 25.7359],
 [0.9177, 0.9177, 0.9535, 0.8798, 0.9181, 26.5012],
 [0.9156, 0.9155, 0.9535, 0.8761, 0.9161, 26.9401]]
```

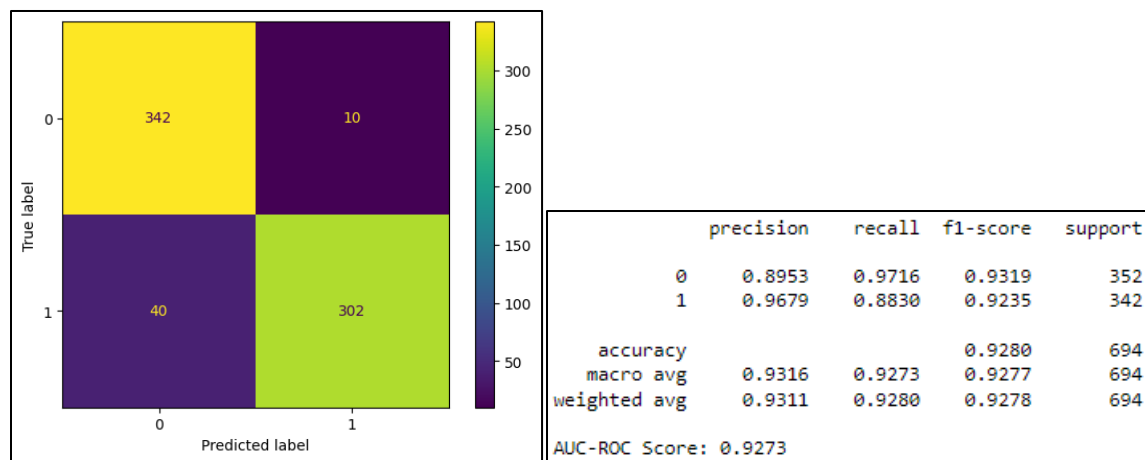
Model Robustness Results:



CatBoost

CatBoost was tested on MalbehavD-V1 to determine implementation on similar datasets (i.e., one with categorical/string data). [Recent simple tests](#) suggests that it is actually capable of handling categorical data, albeit with certain condition(s). Mainly, NaN values shall be converted into a string instead.

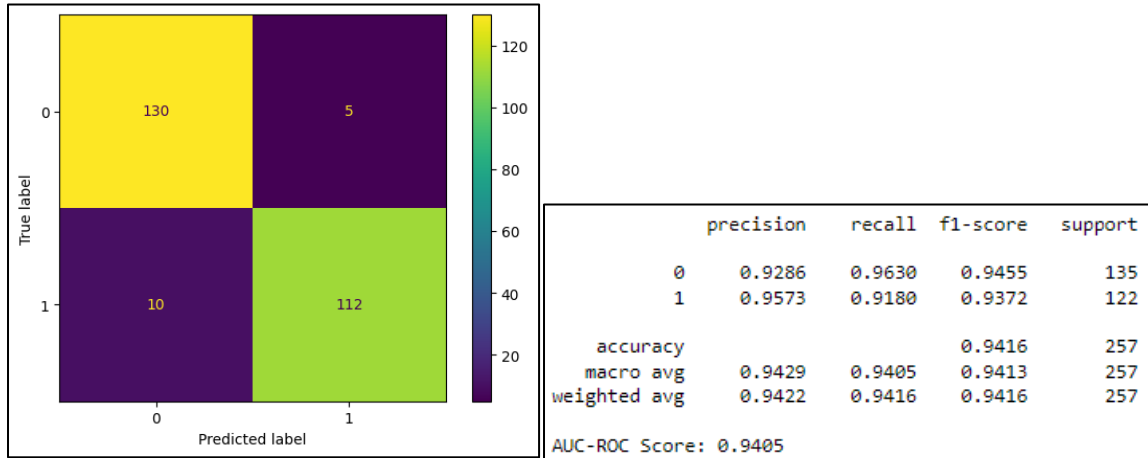
Split Sample Test Results:



K-Folds Test Results:

```
[['accuracy', 'f1_score', 'precision', 'recall', 'roc_auc', 'time'],
 [0.9201, 0.9198, 0.9804, 0.8584, 0.9205, 232.7766],
 [0.9309, 0.9309, 0.9427, 0.9185, 0.931, 224.2083],
 [0.9309, 0.9309, 0.9548, 0.9056, 0.9311, 261.7223],
 [0.9372, 0.9372, 0.9677, 0.9052, 0.9374, 293.2394],
 [0.9156, 0.9155, 0.9447, 0.8836, 0.9157, 238.1865]]
```

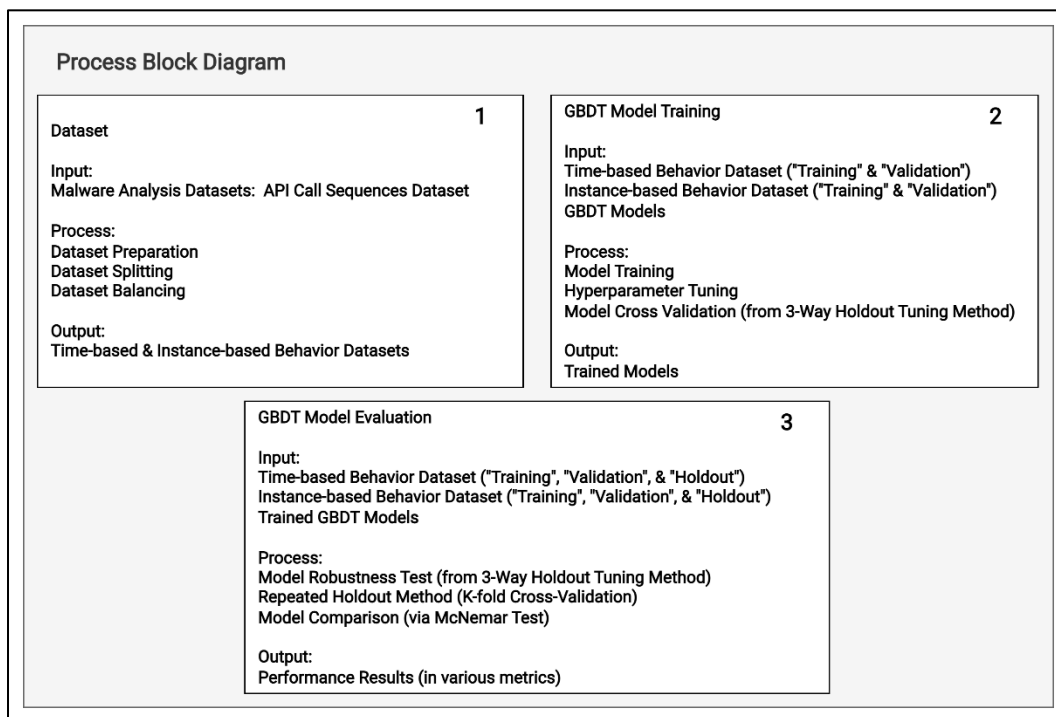
Model Robustness Results:



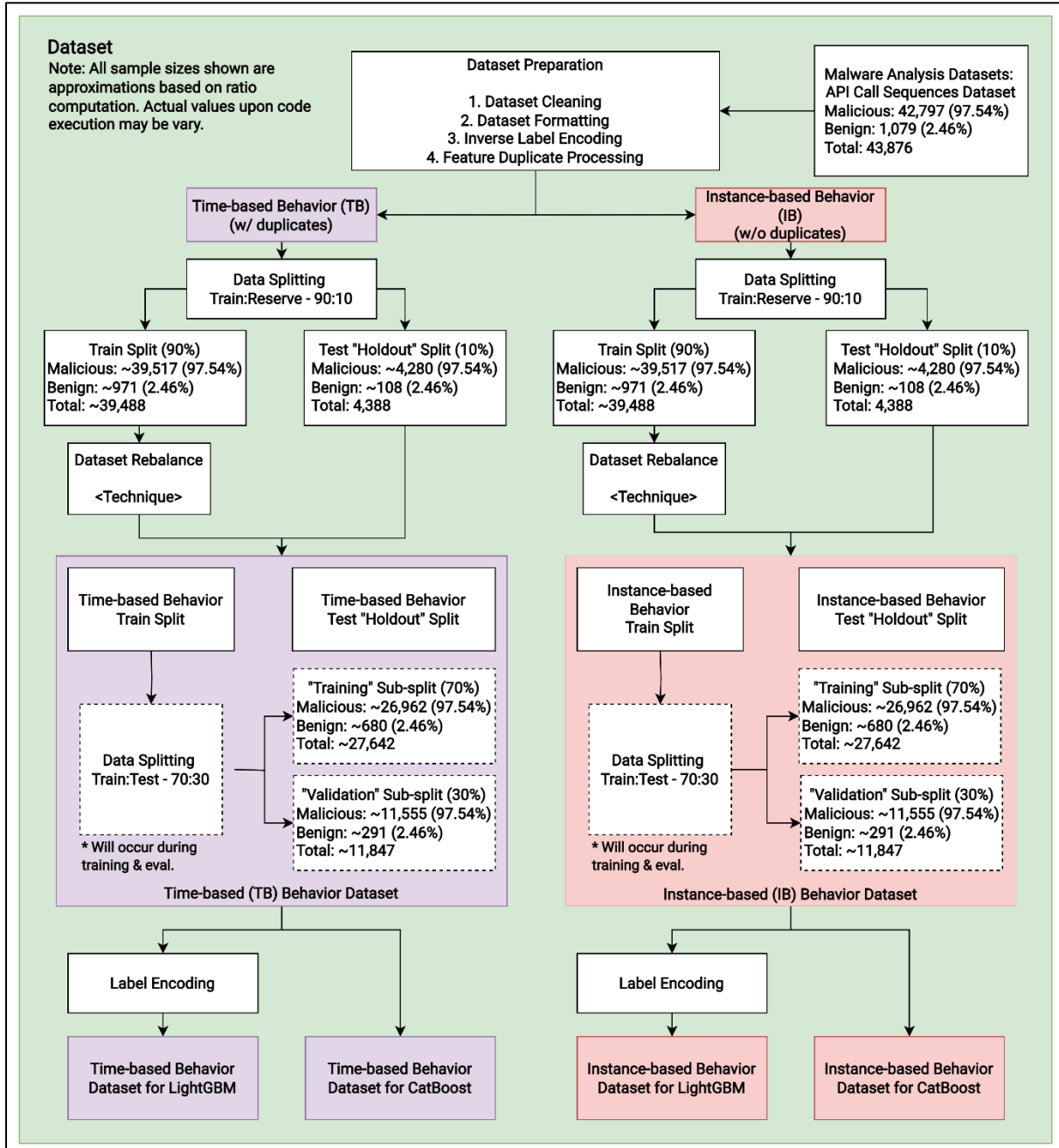
Proposed Process Block Diagram (PBD)

The [proposed Process Block Diagrams](#) will be a total of four diagrams which are the Overview, Dataset, GBDT Model Training, and GBDT Model Evaluation. The last 3 diagrams go in line with the study's Gantt chart where each of the major section of the study is divided into 3 sections for each of the terms (i.e., THES1, THES2, THES3). Note that any of the diagrams shown here are prototypes based on the concepts discussed on CH3 which is then arranged according to the study's objectives, needs, & available resources; & is not yet final. It is also apparent that, in the interest of time and available resources, not all concepts discussed in CH3 will be used in CH4.

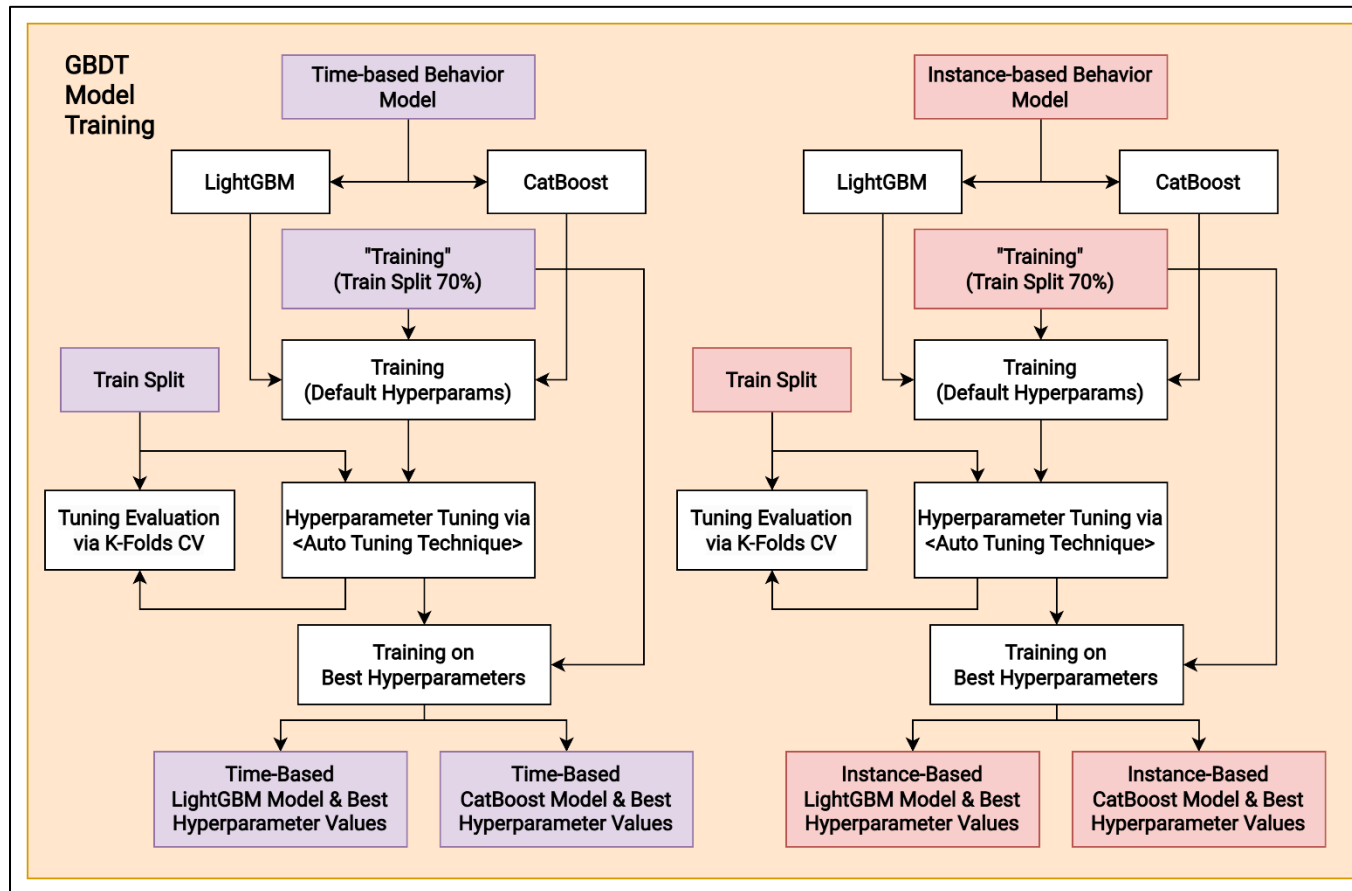
The diagram below shows the Overview PBD where it shows the three major areas of study and the summary of the processes involved in each of the three areas.



The diagram below is about the Dataset. It shows all the processes involved in processing the dataset for it to be useable in the study. The processes shown here are mentioned in CH3. To summarize, the objective of this area of study is to produce the necessary processed datasets which will be used in the next area of the study (i.e., GBDT Model Tuning).



The diagram below is about the GBDT Model Training. It shows all the processes involved in training the model, which is the very essence of the study. The processes shown here are mentioned in CH3. To summarize, the objective of this area of study is to produce the trained models as files to be used on most of the model evaluation (except K-Folds Test) in the next area of the study (i.e., GBDT Model Evaluation).



The diagram below is about the GBDT Model Evaluation. It shows all the processes involved in evaluating the model or the tuning parameters developed in the study. The processes shown here are mentioned in CH3. To summarize, the objective of this area of study is to produce the performance results from various evaluation metrics which can be used to give insights and conclusions as to which model and behavior-type (time-based or instance-based) is better.

