
GTX: Ground Truth eXplanation Dataset

Xiayan Ji Anton Xue Rajeev Alur Oleg Sokolsky Insup Lee Eric Wong

Department of Computer and Information Science

University of Pennsylvania

Philadelphia, PA 19104

{xjiae, antonxue, alur, sokolsky, lee, exwong}@seas.upenn.edu

Abstract

Evaluating the quality of explainability methods is challenging due to the lack of ground truth explanations, and often rely on hand-crafted heuristics. We curate the Ground Truth eXplanation dataset (GTX) to evaluate the alignment of feature attributions with human annotations. These annotations are carefully selected to be directly causal to the ground truth label, which provides an unambiguous goal for human-aligned explainable models. GTX is a diverse benchmark spanning multiple real-world and high dimensional data types (time-series, image, and text). In these settings, the actual explanatory features constitute only a small fraction of the entire feature space. Our analysis finds that common explanation methods overlook the ground truth explanatory features with a worryingly high false negative rate. Our dataset provides a quantitative goal for the future development of feature attribution algorithms: re-aligning explainable models with human explanations. GTX datasets and data loaders publicly available at <https://github.com/xjiae/HDDDS>.

1 Introduction

The size of modern deep networks can easily exceed millions of parameters and hidden units, making it challenging for humans to understand [1]. However, decision makers need to comprehend the model’s reasoning and determine if and when they should rely on these predictions. In higher stakes settings, there are severe repercussions for naively deploying models without fully understanding its reasoning and limitations, such as in diagnosis systems in medicine [2] or legal briefs in judiciaries [3]. To provide some degree of accountability, many post-hoc techniques [4, 5, 6] have been proposed to explain the reasoning behind individual predictions of machine learning models.

One popular class of explanation techniques is feature attributions [7, 8, 9], where given an input, the objective is to assign a score for each feature as it relates to the model’s prediction. Intuitively, the score of a feature is intended to measure the “importance” of said feature towards the model prediction, where larger scores indicate the feature was highly important for making the prediction. Feature attributions have applications in classic machine learning settings such as vision [10, 11], language [12, 13], and reinforcement learning [14, 15], as well as more recent use-cases in industry [16] and law [17].

However, feature attributions rarely come with formal guarantees of behavior [18]. While various metrics have been proposed to evaluate feature attributions, each metric only provides a plausible, partial view into the underlying model’s behavior that need not be accurate. For instance, some metrics progressively remove features with the highest attribution scores and assessing the resulting change [19, 20]. Many of these metrics arose out of necessity because ground-truth explanations for

comparison were unclear or simply not available. As a result, ensuring that feature attributions are of high quality remains a significant challenge despite their widespread usage.

Hence, our work puts forth a measurable and human-aligned target for feature attributions. In order for models to effectively assist humans, it is essential for their decision-making to be aligned with human judgment [21]. Specifically, a feature attribution for an explainable model should identify the ground truth causal features in the *data*.¹ Therefore, we seek to quantify to what degree are models and their feature attributions aligned with human judgment. If explainable models exhibit a high degree of alignment with humans, then their explanations are more *usable* as a proxy for human experts. For example, a doctor could use such a human-aligned attribution to explain a medical diagnoses in lieu of asking a specialist, freeing up the specialist to pursue more challenging tasks.

However, evaluating feature attributions with human judgment is challenging due to the absence of a ground truth explanation [22]. Several benchmarks for evaluating feature attribution methods, such as Captum [23] and OpenXAI [24], have been established using synthetic datasets where the ground truth can be carefully controlled and specified. However, there remains a need to complement these benchmarks with high-dimensional, and real-world datasets that offer a diverse and rich perspective for evaluating of feature attribution algorithms in natural settings. To this end, we have curated a set of real-world datasets that possess ground truth human annotations in domains such as industrial controls [25, 26, 27], artifact evaluation [28], and machine comprehension [29].

Our benchmark, called the Ground Truth eXplanation dataset (GTX), is specifically designed to comprehensively evaluate the human-alignment of feature attribution methods in challenging, real-world settings. Our contribution can be summarized as follows:

- We meticulously clean and process the human annotations to create the Ground Truth Explanation (GTX) dataset. The resulting benchmark spans three prominent data domains: time-series, image, and text, as depicted in Figure 1, but has a standardized and measurable human-alignment goal across all tasks.
- We establish a baseline for the alignment of common feature attribution methods and models with the human annotations using our GTX benchmark.
- In our analysis, we show that existing feature attribution algorithms have a high false negative rate and tend to overlook the true explanatory features. This misalignment highlights the need for future research to achieve more usable explanations.

The remainder of this paper is organized as follows: Section 2 provides a review of existing feature attribution methods and explores relevant XAI benchmarks and datasets. In Section 3, we introduce our dataset, highlighting its unique attributes and characteristics. To showcase potential usage of feature attribution methods on our dataset, we present baseline experiments in Section 4. Subsequently, in Section 5, we discuss the limitations of our dataset and conclude with a summary of our findings.

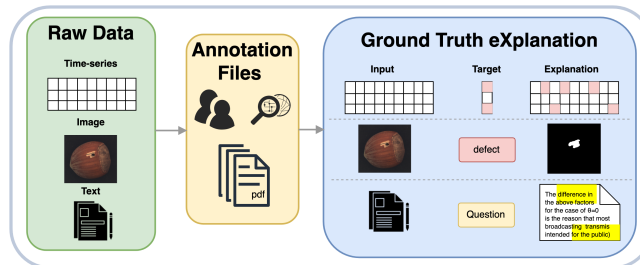


Figure 1: GTX overview. We consolidate raw data of time-series, image and text format. Then, we process the annotation files to obtain the ground truth explanations in column, pixel and clause levels.

¹We note that this criteria (explaining the prediction from patterns in the data) differs from explaining predictions from patterns learned in the model. In the latter, the goal is different—these explanations aim to uncover properties learned in the model which need not be aligned with the ground truth or be usable for humans.

69 2 Related Work

70 2.1 XAI methods on Feature Attribution

71 The existing literature encompasses various taxonomies of Explainable Artificial Intelligence (XAI)
72 methods, each tailored to address specific problems and aspects. In this study, our primary emphasis is
73 on the post-hoc method branch, with a specific focus on feature attribution [5, 7] or feature relevance
74 explanation [6]. Feature attribution refers to the process of determining the importance or contribution
75 of individual features within a dataset or input data to the predictions or output of a machine learning
76 model [7, 8, 9, 10, 30, 31, 32]. Specifically, we reproduce the code of two standard explanation
77 methods for evaluation: Vanilla Gradient [31] and Integrated Gradients [32].

78 2.2 XAI Benchmarks and Datasets

79 Many open-source library implement a handful of feature attribution algorithms, for example, Cap-
80 tum [33] and SHAP bechmarks [34]. However, they do not perform ground-truth based evaluation.
81 Several studies acknowledge the limitation of XAI due to the absence of ground-truth for evaluating
82 explanations [20, 35, 36]. To address this issue, researchers have started introducing ground truth
83 annotations to assess XAI methods. For instance, Amiri et al. [37] propose the use of canonical
84 equations as representations of explanations for evaluating their accuracy. Furthermore, Arras et
85 al. [38] introduce CLEVR-XA, a visual question answering dataset designed specifically for eval-
86 uating neural network explanations in computer vision tasks. OpenXAI [24] offers a transparent
87 evaluation of post hoc model explanations using tabular data and perform faithfulness evaluation with
88 ground truth of the synthetic data. Different from previous works, we present five real-world datasets
89 of different data types, i.e., time-series, image, and text, that can be used to perform evaluation based
90 on human-annotated ground truth explanations.

91 3 Ground Truth eXplanation Dataset

92 In our GTX dataset we consolidate three common types of data in time-series, image, and text,
93 for evaluating feature attribution methods. For time-series data the features correspond to periodic
94 samples of a plant state; for image data the features correspond to pixels of the image; for text data
95 the features correspond to tokens of the text. By analyzing the importance of these features, we can
96 gain insights into the decision-making process of the model from the input level.

97 3.1 Time-series

98 The time-series component of GTX consists of three different datasets from various industrial control
99 settings involving real-world or simulated plants. All three datasets were generated by sampling the
100 plant at a fixed frequency, where for each datapoint the feature values denote either a sensor reading
101 or a controller output.

102 **Hardware-In-the-Loop-based Augmented ICS Security Dataset (HAI)** [25]: The HAI dataset
103 was collected from a realistic industrial control system (ICS) testbed, augmented with a Hardware-
104 In-the-Loop (HIL) simulator for 379.3 hours. The HIL simulator emulates two crucial components
105 of the power generation domain: steam-turbine power generation and pumped-storage hydropower
106 generation, with a total of $m = 86$ features.

107 **Secure Water Treatment Dataset (SWaT)** [26]: The Secure Water Treatment testbed serves as a
108 scaled-down replica of a real-world industrial water treatment plant. It operates at a reduced capacity,
109 producing five gallons per minute of water for over 11 days. The treatment process involves the
110 utilization of membrane-based ultrafiltration and reverse osmosis units for effective water filtration,
111 comprising of $m = 51$ features in total.

112 **Water Distribution Dataset (WADI)** [27]: WADI is an extension of the SWaT testbed featuring
 113 additional components and functionalities such as chemical dosing systems, booster pumps and valves,
 114 as well as instrumentation and analyzers. It is collected over 16 days with $m = 127$ dimensions.

115 At various time points a cyber attack (e.g. altering sensor readings) or a physical attack (e.g. altering
 116 water flow) is performed, which allows one to obtain a ground truth of whether a plant state is to be
 117 considered “normal” or “attacked”. The attacks target specific sets of equipment and have precise
 118 start and end times, allowing us to obtain accurate ground truth explanations for the features during
 119 the attacks. We manually process annotation files in both PDF and Excel formats, which contain the
 120 start and end times of each attack. These annotations are then aligned with the timestamps of the raw
 121 data records.

122 Each dataset is a sequence of periodic samples $(x_1, y_1), \dots, (x_T, y_T)$ where at each time step
 123 $t = 1, \dots, T$ the observations $x_t \in \mathbb{R}^m$ denotes plant state while the label $y_t \in \{0, 1\}$ is an indicator
 124 of whether the plant was attacked ($y_t = 1$) or if the behavior is normal ($y_t = 0$). We use a binary
 125 mask $a_t \in \{0, 1\}^m$ as the ground truth explanation to denote which input feature is explanatory of
 126 an attack. If $y_t = 1$, then $(a_t)^i = 1$ implies that feature i is involved in attack at time t — which
 127 we know from the annotations supplied with the original datasets. If $y_t = 0$, then we write $a_t = 0$,
 128 the zeros vector, to mean that no attack occurred. In summary, our dataset loader provides three key
 129 objects at each time step: the plant features x_t , the attack indicator y_t , and the explanation a_t .

130 To facilitate the use of our dataset with machine learning models in PyTorch [39], we wrap the raw
 131 data using the `torch.utils.data.Dataset` class. Below is an example code snippet demonstrating
 132 its usage with the HAI dataset:

```
133 bundle = get_data_bundle("hai", window_size=100, train_batch_size=32)
134 train_dataloader = bundle["train_dataloader"]
135 x, y, a = next(iter(train_dataloader))
136 # x.shape==(32,100,86), y.shape==(32,), a.shape==(32,100,86)
137
```

139 3.2 Image

140 MVTec-AD [28] is an industrial inspection dataset designed for benchmarking defects detection
 141 methods. It consists of a 15 categories with a total of more than 5000 high-resolution ($3 \times 1024 \times 1024$)
 142 images. Each category includes a set of defect-free training images and a test set containing images
 143 with different types of defects, as well as defect-free images. The dataset provides pixel-accurate
 144 ground truth annotations for the defect regions, which have been carefully annotated and reviewed by
 145 the authors to align with human interpretation of real-world defects.

146 We allow the user to specify an input size $d \leq 1024$ to down-sample an image to $m = 3 \times d \times d$
 147 features, i.e. 3 color channels with a side-length of d pixels. Each image is correspondingly labeled
 148 with whether it has a defect ($y = 1$) or not ($y = 0$). The ground truth explanation is a bitmask
 149 $a \in \{0, 1\}^{d \times d}$ denoting which positions are defects; if $(a)^{ij} = 1$, then this means that the pixel at
 150 position (i, j) is part of the defect. If $y = 0$, then $a = 0$, the zeros matrix, indicating no defects. The
 151 objects returned by the dataset are the down-sampled image x , the defect label y , and the ground
 152 truth explanation a .

153 As there are 15 image categories in total, MVTec is in fact a collection of 15 different datasets. We
 154 implement the MVTec dataset with `torch.utils.data.Dataset`, and showcase its use below.

```
155 bundle = get_data_bundle("mvtec", input_size=256, train_batch_size=32)
156 train_dataloader = bundle["train_dataloader"]
157 x, y, a = next(iter(train_dataloader))
158 # x.shape==(32,3,256,256), y.shape==(32,), a.shape==(32,1,256,256)
159
```

161 Here hazelnut is one of 15 admissible classes among:

```
162         bottle, cable, capsule, carpet, grid, hazelnut, leather, metal_nut,
163         pill, screw, tile, toothbrush, transistor, wood, zipper
```

Here `input_size` is the dimension d to which we downsample, and the last flag of `is_train=True` selects only images that are non-defect; if `is_train=False` then the selection is mixed.

3.3 Text

SQuAD (Stanford Question Answering Dataset) [29] is a widely used reading comprehension dataset that includes 107,785 question-answer pairs based on 536 Wikipedia articles. The dataset was generated by crowdworkers who formulated questions and provided specific text segments or spans as answers. The answers have undergone rigorous crowdworkers selection, additional answer collection, and manual crosscheck processes, making them reliable ground truth explanations for the corresponding questions.

The questions are concatenated with a context, such that model inputs have the form $x = (x_q, x_c)$, where x_q are the question tokens and x_c are the context tokens. The output of a question-answering model is to identify a range of indices to highlight in the conjoined input x that constitutes as the answer. As such, the output of a model on SQuAD is not a binary value as in the time-series and image data, but instead a start-index and an end-index that denotes which tokens to highlight. For a particular x in the dataset, the ground truth then consists of a pairing $a = (a_s, a_e)$, where a_s, a_e are integers that denote the highlight start and end indices, respectively, for x .

We implemented the SQuAD dataset with `torch.utils.data.Dataset` as follows, where we demonstrate tokenization with the RoBERTa [40] base tokenizer to sequence lengths of 384.

```
# Use the "roberta-base" tokenizer from Hugging Face
bundle = get_data_bundle("squad", tokenizer_or_name="roberta-base",
                        train_batch_size=32)
train_dataloader = bundle["train_dataloader"]
input_ids, attn_mask, token_type_ids, start_pos, end_pos = next(iter(
    train_dataloader))
# input_ids.shape==(32,384), the default token sequence length
```

Each item within a SQuAD dataset contains a number of information relevant for a language-model transformer, among them: `item[0]` corresponds to x , which we emphasize is the concatenation of the question tokens and the context tokens; `item[3]` and `item[4]` correspond to the start and end position indices, respectively. We use the defaults supplied with `tensorflow_datasets` and `transformers` to determine the train-test split.

3.4 Dataset Statistics

In Table 1, we provide a summary of key statistics pertaining to the datasets. This includes information on the feature dimensions, the number of positive instances representing attacks or defects, the number of negative instances, and the corresponding positive ratio. Notably, the class distribution exhibits an imbalance, signifying a discrepancy in the distribution between negative and positive instances. We use “positive” to denote data for which $y = 1$, and “negative” to denote data for which $y = 0$.

However, it is important to note that the SQuAD dataset presents unique characteristics that distinguish it from the other datasets. The variable lengths of paragraphs contribute to the variability in feature dimensions. Additionally, the SQuAD dataset does not differentiate between positive and negative instances, making it unsuitable for inclusion in Table 1 for comparative purposes.

Data	Features Dimension	Positive Count	Negative Count	Total Count	Positive Ratio
HAI	86	12,030	1,353,572	1,365,602	0.88%
SWaT	51	54,621	892,098	946,719	5.77%
WADI	127	5,134	1,377,268	1,382,402	0.37%
MVTec	$3d^2$	1256	4094	5350	23.48%

Table 1: Basic statistics of HAI, SWaT, WADI, and MVTec. When $y = 1$ we say that the datapoint is positive; when $y = 0$ we say that it is negative.

Furthermore, we present several key statistics of the human annotations within our dataset. Specifically, we report the total number of human annotations conducted, the average number of explanatory features per input, and the ratio of explanatory features to the entire feature space. This ratio serves as a crucial metric for assessing the class imbalance between explanatory and non-explanatory features in the annotations. For the SQuAD dataset, we gather and estimate the summary statistics based on information provided in their papers [29, 41].

The presence of class imbalance poses significant challenges for feature attribution methods, as they aim to accurately identify and attribute the importance of each feature in the prediction process. When the number of explanatory features is significantly lower than the number of non-explanatory features, it can lead to biased attributions and potentially misleading interpretations of the model’s behavior.

Data	Annotation Count	Average Count	Explanatory Ratio
HAI	1,034,580	1.00	1.17%
SWaT	2,785,671	1.07	2.10%
WADI	652,018	1.93	1.52%
MVTec	1,317,011,456	45950.49	4.38%
SQuAD	107,785	4.64	3.10%

Table 2: Feature statistics of HAI, SWaT, WADI, MVTEC and SQuAD.

3.5 Task Definition

We formulate the task as predicting which features of an input $x \in \mathcal{X}, \mathcal{X} \subseteq \mathbb{R}^m$ related to the target $y \in \mathcal{Y}$. Specifically, a feature attribution model $A : \mathcal{X} \rightarrow [0, 1]^m$ maps an input x to an m -dimensional vector $\hat{a} = \hat{A}(x) \in [0, 1]^m$, where each element is a score representing the degree that the corresponding feature is explanatory of y . For each input x and target y , our dataset has a m -bit vector that encodes the ground truth annotation function $A : \mathcal{X} \rightarrow \{0, 1\}^m$. It maps each input x to the human annotation $a = A(x) \in \{0, 1\}^m$. By comparing \hat{a} and a , we can directly evaluate the performance of the feature attribution model.

4 Experiments

4.1 Predictive models

In our experiments involving time-series data, we utilize the standard implementation of well-established logistic regression model (LR) [42] and Long Short-term Memory networks (LSTM). In our analysis of the MVTEC dataset, we utilize Fastflow [43], a CNN that employs ResNet18 as its backbone for image feature extraction. The SQuAD dataset is processed using the widely-used RoBERTa model [40]. These models are chosen to facilitate the application of diverse feature attribution methods.

4.2 Feature Attribution methods

In our study, we employ several widely used feature attribution techniques. Specifically, we apply the vanilla gradient (GRAD) [31] and integrated gradient (INTG) [32] to evaluate the feature attribution performance. The objective of this analysis is to assess the importance of individual features in the decision-making process of the machine learning models. To quantify the effectiveness of the feature attribution techniques in a binary manner, we establish a threshold on the feature attribution scores. This threshold is determined by maximizing the F1-score, a widely utilized metric that balances precision and recall.

4.3 Training

In our time-series data study, we employed two distinct training strategies. Firstly, we randomly divided 70% of the dataset as the training set and train the Logistic Regression (LR) model. Secondly,

we preserved the temporal order by training a three-layered Long Short-Term Memory (LSTM) model on the sliding window version of the dataset. The training ratio of 70% was maintained, and a window size of 100 with a stride of one was used. For the image data, the FastFlow model was trained on 70% of the MVTEC dataset. We sampled both negative and positive instances for the time-series and image training sets. The text data was trained using the RoBERTa model on the default training set, which accounted for 82% of the SQuAD dataset.

All models were trained for five epochs with a learning rate of 10^{-6} . The model with the best validation accuracy was selected. All experiments are run with 4 Nvidia 2080Ti GPU, 80 vCPUs, a processor Intel(R) % Xeon(R) Gold 6148 @ 2.4 GHz and 768GiB of RAM. Further details on the model architecture can be found in the Appendix. It is important to note that hyperparameter tuning was not the main focus of our study, as our primary objective was to showcase the utility of our dataset.

4.4 Metrics

By employing thresholds that optimize the F1-score on feature attribution, we obtain binary predictions for individual features. In this context, we consider a prediction of ‘1’ a positive outcome (explanatory), while ‘0’ denotes a negative outcome (not explanatory). Let $a = A(x) = \{a^i \mid i = 1, \dots, m\} \in \{0, 1\}^m$ represent the ground truth annotation for $x \in \mathbb{R}^m$. Let $\hat{a} = \hat{A}(x) = \{\hat{a}^i \mid i = 1, \dots, m\} \in \{0, 1\}^m$ denote the prediction generated by the attribution models and the threshold. We then concatenate all the features of all the inputs and compute several evaluation metrics, including False Positive Rate (FPR), False Negative Rate (FNR), Accuracy (ACC), and F1-score. The computation formulas for these metrics are presented below:

$$\text{FPR} = \frac{|\{\hat{a}^i = 1 \mid a^i = 0\}|}{|\{a^i = 0\}|}, \text{FNR} = \frac{|\{\hat{a}^i = 0 \mid a^i = 1\}|}{|\{a^i = 1\}|} \quad (1)$$

$$\text{ACC} = \frac{|\{\hat{a}^i = 0 \mid a^i = 0\}| + |\{\hat{a}^i = 1 \mid a^i = 1\}|}{|\{a^i = 0\}| + |\{a^i = 1\}|} \quad (2)$$

$$\text{F1-score} = \frac{2|\{\hat{a}^i = 1 \mid a^i = 1\}|}{2|\{\hat{a}^i = 1 \mid a^i = 1\}| + |\{\hat{a}^i = 1 \mid a^i = 0\}| + |\{\hat{a}^i = 0 \mid a^i = 1\}|}. \quad (3)$$

4.5 Results

We randomly sampled 100 instances on individual datasets and conducted the experiments using 20 random seeds. The obtained results were then analyzed by reporting the mean and standard error metric.

Time-series The results obtained from our experiments on time-series data have revealed that GRAD and INTG miss-classify the explanatory feature as non-explanatory, which leads to a higher FNR (+42.23%) than FPR on average, as shown in Table 3, 4, 5. In addition, we observe that GRAD has a better performance in general, with higher accuracy (+17.71%) and F1-score (+11.10%) performance than that of the INTG. This presents a serious problem as it undermines the reliability

Models	Attribution	FPR	FNR	ACC	F1-score
LR	GRAD	0.02 ± 0.02	0.89 ± 0.09	0.97 ± 0.02	0.98 ± 0.01
	INTG	0.49 ± 0.03	0.99 ± 0.01	0.51 ± 0.03	0.67 ± 0.03
LSTM	GRAD	0.06 ± 0.10	0.87 ± 0.19	0.94 ± 0.10	0.96 ± 0.06
	INTG	0.01 ± 0.00	0.89 ± 0.02	0.98 ± 0.00	0.99 ± 0.00

Table 3: Results for HAI dataset.

and effectiveness of the attribution methods in correctly identifying the features that contribute to the model’s decision-making process. Consequently, it highlights the need for further improvement and development of attribution algorithms to address this challenge and enhance their capability to accurately identify and attribute the explanatory features.

Models	Attribution	FPR	FNR	ACC	F1-score
LR	GRAD	0.03 \pm 0.00	0.50 \pm 0.05	0.96 \pm 0.00	0.97 \pm 0.00
	INTG	0.52 \pm 0.03	0.90 \pm 0.03	0.48 \pm 0.03	0.64 \pm 0.03
LSTM	GRAD	0.31 \pm 0.36	0.66 \pm 0.35	0.69 \pm 0.36	0.73 \pm 0.37
	INTG	0.53 \pm 0.01	0.66 \pm 0.04	0.47 \pm 0.01	0.63 \pm 0.01

Table 4: Results for SWaT dataset.

Models	Attribution	FPR	FNR	ACC	F1-score
LR	GRAD	0.58 \pm 0.07	0.28 \pm 0.07	0.42 \pm 0.07	0.58 \pm 0.07
	INTG	0.50 \pm 0.03	0.94 \pm 0.03	0.49 \pm 0.03	0.66 \pm 0.03
LSTM	GRAD	0.31 \pm 0.11	0.46 \pm 0.14	0.69 \pm 0.11	0.80 \pm 0.09
	INTG	0.33 \pm 0.24	0.71 \pm 0.20	0.67 \pm 0.24	0.77 \pm 0.16

Table 5: Results for WADI dataset.

Image For the image data, Table 11 demonstrates that a higher FNR (+69.97%) than FPR is also observed, suggesting that the attribution methods fail to capture all the explanatory pixels. Different from the time-series dataset, INTG has a better performance than GRAD method, with a higher accuracy (+3.29%) and F1-score (+2.15%).

Models	Attribution	FPR	FNR	ACC	F1-score
FastFlow	GRAD	0.14 \pm 0.11	0.79 \pm 0.11	0.86 \pm 0.10	0.91 \pm 0.07
	INTG	0.10 \pm 0.05	0.85 \pm 0.06	0.89 \pm 0.05	0.93 \pm 0.03

Table 6: Results for MVTec dataset.

283

Text The result of the SQuAD dataset can be found in Table 12. As with previous datasets, a notable observation is the presence of a high FNR compared with the image dataset. The same result for INTG and GRAD could be due to the similar gradient computations by the RoBERTa model. However, the results may vary depending on the predictive model architecture, and the complexity of the explanation task in practice.

Overall, the GRAD method demonstrates a slightly better performance than INTG, exhibiting a higher accuracy (+9.97%) and F1-score (+6.23%) on average. In Figure 2, we present a comparative analysis of average FPRs and FNRs for different attribution methods and datasets. The figure highlights that SQuAD exhibits the highest FNR while WADI showcases the highest FPR on average.

More figures and results on other attribution methods (e.g. SHAP [7] and LIME [8]) can be found in the Appendix. Our experiments are illustrative in nature, running with different machines or configurations may yield slightly different results. However, the overall trends and patterns observed in the data should remain similar and consistent.

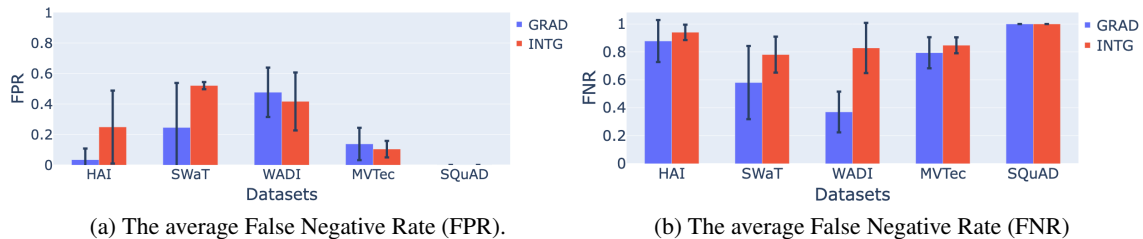


Figure 2: Average error rates comparison for INTG and GRAD across all datasets.

Models	Attribution	FPR	FNR	ACC	F1-score
RoBERTa	GRAD	0.00 ± 0.00	1.00 ± 0.00	0.87 ± 0.00	0.81 ± 0.00
	INTG	0.00 ± 0.00	1.00 ± 0.00	0.87 ± 0.00	0.81 ± 0.00

Table 7: Results for SQuAD dataset.

5 Conclusion

One limitation of our dataset is the absence of real-world graph data. However, we are actively searching and we will update our repository once we find suitable datasets. In summary, our GTX dataset includes time-series, image, and text data, along with detailed feature-wise ground truth explanations. We have established a baseline for aligning common feature attribution algorithms with human annotation of the actual explanatory features, which takes up a relatively small proportion in real-world datasets. Our experiments have revealed a significant challenge posed by a higher FNR than FPR in existing feature attribution methods, emphasizing the need for improvements to accurately identify the true explanatory features. With its comprehensive collection and diverse data types, our dataset is a valuable resource for the XAI community, facilitating quantitative evaluation and advancements in feature attribution algorithms.

References

- [1] Zhengxuan Wu, Atticus Geiger, Christopher Potts, and Noah D Goodman. Interpretability at scale: Identifying causal mechanisms in alpaca. *arXiv preprint arXiv:2305.08809*, 2023.
- [2] Julia Amann, Alessandro Blasimme, Effy Vayena, Dietmar Frey, Vince I Madai, and Precise4Q Consortium. Explainability for artificial intelligence in healthcare: a multidisciplinary perspective. *BMC medical informatics and decision making*, 20:1–9, 2020.
- [3] Ashley Deeks. The judicial demand for explainable artificial intelligence. *Columbia Law Review*, 119(7):1829–1850, 2019.
- [4] W James Murdoch, Chandan Singh, Karl Kumbier, Reza Abbasi-Asl, and Bin Yu. Definitions, methods, and applications in interpretable machine learning. *Proceedings of the National Academy of Sciences*, 116(44):22071–22080, 2019.
- [5] Pantelis Linardatos, Vasilis Papastefanopoulos, and Sotiris Kotsiantis. Explainable ai: A review of machine learning interpretability methods. *Entropy*, 23(1):18, 2020.
- [6] Alejandro Barredo Arrieta, Natalia Díaz-Rodríguez, Javier Del Ser, Adrien Bannetot, Siham Tabik, Alberto Barbado, Salvador García, Sergio Gil-López, Daniel Molina, Richard Benjamins, et al. Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai. *Information fusion*, 58:82–115, 2020.
- [7] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *Advances in neural information processing systems*, 30, 2017.
- [8] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should i trust you?": Explaining the predictions of any classifier, 2016.
- [9] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Anchors: High-precision model-agnostic explanations. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- [10] Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. Smoothgrad: removing noise by adding noise, 2017.

- [11] Vanessa Buhrmester, David Münch, and Michael Arens. Analysis of explainers of black box deep neural networks for computer vision: A survey. *Machine Learning and Knowledge Extraction*, 3(4):966–989, 2021.
- [12] Yongfeng Zhang, Guokun Lai, Min Zhang, Yi Zhang, Yiqun Liu, and Shaoping Ma. Explicit factor models for explainable recommendation based on phrase-level sentiment analysis. In *Proceedings of the 37th international ACM SIGIR conference on Research & development in information retrieval*, pages 83–92, 2014.
- [13] Marina Danilevsky, Kun Qian, Ranit Aharonov, Yannis Katsis, Ban Kawas, and Prithviraj Sen. A survey of the state of explainable ai for natural language processing. *arXiv preprint arXiv:2010.00711*, 2020.
- [14] Nikaash Puri, Sukriti Verma, Piyush Gupta, Dhruv Kayastha, Shripad Deshmukh, Balaji Krishnamurthy, and Sameer Singh. Explain your move: Understanding agent actions using specific and relevant feature attribution. *arXiv preprint arXiv:1912.12191*, 2019.
- [15] Lei He, Nabil Aouf, and Bifeng Song. Explainable deep reinforcement learning for uav autonomous path planning. *Aerospace science and technology*, 118:107052, 2021.
- [16] Wenqian Liu, Runze Li, Meng Zheng, Srikrishna Karanam, Ziyang Wu, Bir Bhanu, Richard J. Radke, and Octavia Camps. Towards visually explaining variational autoencoders. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [17] Katie Atkinson, Trevor Bench-Capon, and Danushka Bollegala. Explanation in ai and law: Past, present and future. *Artificial Intelligence*, 289:103387, 2020.
- [18] Meike Nauta, Jan Trienes, Shreyasi Pathak, Elisa Nguyen, Michelle Peters, Yasmin Schmitt, Jörg Schlötterer, Maurice van Keulen, and Christin Seifert. From anecdotal evidence to quantitative evaluation methods: A systematic review on evaluating explainable ai. *arXiv preprint arXiv:2201.08164*, 2022.
- [19] Aya Abdelsalam Ismail, Mohamed Gunady, Hector Corrada Bravo, and Soheil Feizi. Benchmarking deep learning interpretability in time series predictions. *Advances in neural information processing systems*, 33:6441–6452, 2020.
- [20] Sara Hooker, Dumitru Erhan, Pieter-Jan Kindermans, and Been Kim. A benchmark for interpretability methods in deep neural networks. *Advances in neural information processing systems*, 32, 2019.
- [21] Giang Nguyen, Daeyoung Kim, and Anh Nguyen. The effectiveness of feature attribution methods and its correlation with automatic evaluation scores. *Advances in Neural Information Processing Systems*, 34:26422–26436, 2021.
- [22] Yilun Zhou, Serena Booth, Marco Tulio Ribeiro, and Julie Shah. Do feature attribution methods correctly attribute features? In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 9623–9633, 2022.
- [23] Narine Kokhlikyan, Vivek Miglani, Miguel Martin, Edward Wang, Bilal Alsallakh, Jonathan Reynolds, Alexander Melnikov, Natalia Kliushkina, Carlos Araya, Siqi Yan, et al. Captum: A unified and generic model interpretability library for pytorch. *arXiv preprint arXiv:2009.07896*, 2020.
- [24] Chirag Agarwal, Satyapriya Krishna, Eshika Saxena, Martin Pawelczyk, Nari Johnson, Isha Puri, Marinka Zitnik, and Himabindu Lakkaraju. Openxai: Towards a transparent evaluation of model explanations. *Advances in Neural Information Processing Systems*, 35:15784–15799, 2022.

- [25] Hyeok-Ki Shin, Woomyo Lee, Jeong-Han Yun, and Byung-Gil Min. Two ics security datasets and anomaly detection contest on the hil-based augmented ics testbed. In *Cyber Security Experimentation and Test Workshop*, CSET '21, page 36–40, New York, NY, USA, 2021. Association for Computing Machinery.
- [26] Aditya P Mathur and Nils Ole Tippenhauer. Swat: A water treatment testbed for research and training on ics security. In *2016 international workshop on cyber-physical systems for smart water networks (CySWater)*, pages 31–36. IEEE, 2016.
- [27] Chuadhry Mujeeb Ahmed, Venkata Reddy Palleti, and Aditya P Mathur. Wadi: a water distribution testbed for research in the design of secure cyber physical systems. In *Proceedings of the 3rd international workshop on cyber-physical systems for smart water networks*, pages 25–28, 2017.
- [28] Paul Bergmann, Michael Fauser, David Sattlegger, and Carsten Steger. Mvtec ad—a comprehensive real-world dataset for unsupervised anomaly detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9592–9600, 2019.
- [29] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ Questions for Machine Comprehension of Text. *arXiv e-prints*, page arXiv:1606.05250, 2016.
- [30] Avanti Shrikumar, Peyton Greenside, Anna Shcherbina, and Anshul Kundaje. Not just a black box: Learning important features through propagating activation differences, 2017.
- [31] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*, 2013.
- [32] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *International conference on machine learning*, pages 3319–3328. PMLR, 2017.
- [33] Narine Kokhlikyan, Vivek Miglani, Miguel Martin, Edward Wang, Bilal Alsallakh, Jonathan Reynolds, Alexander Melnikov, Natalia Kliushkina, Carlos Araya, Siqi Yan, and Orion Reblitz-Richardson. Captum: A unified and generic model interpretability library for pytorch, 2020.
- [34] Shap benchmark. <https://shap.readthedocs.io/en/latest/index.html>. [Accessed 17-May-2023].
- [35] Finale Doshi-Velez and Been Kim. Towards a rigorous science of interpretable machine learning, 2017.
- [36] Richard Tomsett, Dan Harborne, Supriyo Chakraborty, Prudhvi Gurram, and Alun Preece. Sanity checks for saliency metrics, 2019.
- [37] Shideh Shams Amiri, Rosina O. Weber, Prateek Goel, Owen Brooks, Archer Gandle, Brian Kitchell, and Aaron Zehm. Data representing ground-truth explanations to evaluate xai methods, 2020.
- [38] Leila Arras, Ahmed Osman, and Wojciech Samek. Clevr-xai: a benchmark dataset for the ground truth evaluation of neural network explanations. *Information Fusion*, 81:14–40, 2022.
- [39] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- [40] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach, 2019.

- [41] Jason Freeman and Raine Hoover. Question answering with squad : Variations on multi-perspective context matching. 2017.
- [42] sklearn logistic regression. https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html.
- [43] Jiawei Yu, Ye Zheng, Xiang Wang, Wei Li, Yushuang Wu, Rui Zhao, and Liwei Wu. Fastflow: Unsupervised anomaly detection and localization via 2d normalizing flows, 2021.
- [44] Dan Hendrycks, Collin Burns, Anya Chen, and Spencer Ball. Cuad: An expert-annotated nlp dataset for legal contract review. *arXiv preprint arXiv:2103.06268*, 2021.

A Appendix

A.1 Model Architecture

Three-layered LSTM model

```
SimpleLSTM(
(lstm1): LSTM(num_features, 128)
(lstm2): LSTM(128, 128)
(lstm3): LSTM(128, 128)
(linear): Linear(128, 2))
```

We used the FastFlow with ResNet18 backbone and RoBERTa-base model. For details of model configuration, please refer to their papers [43, 40].

A.2 More Experiment Results

In this section, we present additional experimental results for the LIME method [8] and the SHAP method [7] applied to all the datasets, utilizing their respective predictive models. It is worth noting that the issue persists in both methods, whereby the false negative rate (FNR) remains greater than the false positive rate (FPR).

Models	Attribution	FPR	FNR	ACC	F1-score
LSTM	LIME	0.19 ± 0.13	0.74 ± 0.13	0.81 ± 0.13	0.88 ± 0.09
	SHAP	0.53 ± 0.04	0.48 ± 0.03	0.47 ± 0.04	0.64 ± 0.03

Table 8: Results on LIME and SHAP for HAI dataset.

Models	Attribution	FPR	FNR	ACC	F1-score
LSTM	LIME	0.37 ± 0.20	0.63 ± 0.19	0.63 ± 0.20	0.75 ± 0.14
	SHAP	0.53 ± 0.02	0.51 ± 0.02	0.47 ± 0.02	0.63 ± 0.02

Table 9: Results on LIME and SHAP for SWaT dataset.

Models	Attribution	FPR	FNR	ACC	F1-score
LSTM	LIME	0.44 ± 0.15	0.55 ± 0.15	0.56 ± 0.15	0.70 ± 0.10
	SHAP	0.18 ± 0.16	0.79 ± 0.15	0.82 ± 0.16	0.88 ± 0.11

Table 10: Results on LIME and SHAP for WADI dataset.

Models	Attribution	FPR	FNR	ACC	F1-score
FastFlow	LIME	0.28 ± 0.25	0.57 ± 0.23	0.72 ± 0.25	0.80 ± 0.20
	SHAP	0.18 ± 0.21	0.73 ± 0.18	0.82 ± 0.21	0.87 ± 0.18

Table 11: Results on LIME and SHAP for MVTec dataset.

Models	Attribution	FPR	FNR	ACC	F1-score
RoBERTa	LIME	0.01 ± 0.03	1.00 ± 0.02	0.99 ± 0.03	0.98 ± 0.01
	SHAP	0.00 ± 0.00	1.00 ± 0.00	0.99 ± 0.00	0.99 ± 0.00

Table 12: Results on LIME and SHAP for SQuAD dataset.

Users are free to apply our dataset to modern models such as the tabular version of Transformer, Vision Transformers, and GPT-2. Since we provide PyTorch compatibility, running these models is straightforward.

A.3 URL to website

URL: <https://github.com/xjiae/HDDDS>.

A.4 Author statement

We bear all responsibility in case of violation of rights, and confirm that we will use the MIT License.

A.5 Dataset documentation, intended use and metadata

We adopt the framework data cards, please find the requested information below. For better visual display, please visit this link: <https://github.com/xjiae/HDDDS/blob/main/description.md>.

The Ground Truth eXplanation (GTX) dataset is a curated collection that addresses the challenge of evaluating the quality of explainability methods. Existing approaches often lack ground truth explanations and heavily rely on hand-crafted heuristics. In response, the GTX dataset has been created to assess the alignment of feature attributions with human annotations. It contains time-series data (HAI, SWaT, WADI) from the industrial control domain, image data (MVTec) from the defect inspection domain, and text data (SQuAD) from the machine comprehension domain.

Dataset Link Dataset Link: HAI, SWaT, WADI, MVTec, SQuAD.

Data Card Author(s)

- **Xiayan Ji, University of Pennsylvania:** (Manager)
- **Anton Xue, University of Pennsylvania:** (Manager)

A.5.1 Authorship

Dataset Owners

Team(s) University of Pennsylvania

Author(s)

- Xiayan Ji, Ph.D. Student, University of Pennsylvania, 2023
- Anton Xue, Ph.D. Student, University of Pennsylvania, 2023
- Rajeev Alur, Professor, University of Pennsylvania, 2023
- Oleg Sokolsky, Professor, University of Pennsylvania, 2023
- Insup Lee, Professor, University of Pennsylvania, 2023
- Eric Wong, Assistant Professor, University of Pennsylvania, 2023

477 **Contact Detail(s)**

- 478 • **Point of Contact:** Xiayan Ji
479 • **Affiliation:** University of Pennsylvania
480 • **Contact:** xjiae@seas.upenn.edu

481 **A.5.2 Dataset Overview**

- 482 • Data about places and objects
483 • Synthetically generated data
484 • Data about systems or products and their behaviors

485

486 **Dataset Snapshot**

Category	Data
Size of Dataset	12 GB
Number of Instances	3,798,242
Number of Labels (explanation)	5,951,278,880
Average Labels Per Instance	1566.85
Algorithmic Labels	4,629,687,370
Human Labels	1,321,591,510

487 **Dataset Summary:** time-series, image and text data with ground truth explanation labels.

488 **Content Description** Each content contains an input data (x), a target label (y) and an explanation
489 (a).

490 **Additional Notes:** for SQuAD, the format is slightly different, the input and target are combined
491 together to better be fitted to a language model. In addition, the explanation is in the form of a start
492 and end position.

493 **Risk Type(s)**

- 494 • No Known Risks

495 **A.5.3 Dataset Version and Maintenance**

496 **Maintenance Status** **Regularly Updated** - New versions of the dataset have been or will continue
497 to be made available.

498 **Version Details** **Current Version:** 1.0

499 **Last Updated:** 06/2023

500 **Release Date:** 06/2023

501 **Maintenance Plan** In our maintenance plan, our primary focus will be on preserving and leveraging
502 the existing data that we have collected. This involves ensuring the integrity and security of the data
503 through regular backups, implementing robust data storage practices, and conducting periodic audits
504 to identify any potential issues or anomalies. Additionally, we recognize the growing importance
505 of graph datasets in various domains. To capitalize on this, we will actively explore and evaluate
506 potential graph datasets that align with our needs and objectives. This includes seeking out reliable

sources, assessing the quality and relevance of the data, and integrating suitable graph datasets into our existing infrastructure. By incorporating graph datasets, we aim to enhance the depth and breadth of our analysis, uncover hidden patterns and relationships, and gain valuable insights that can drive informed decision-making and optimize our operations. In addition, we are aware that the SQuAD dataset does not have a clear classification task and may not align well with the remaining dataset. We are also exploring the Contract Understanding Atticus Dataset (CUAD) [44] to see if we can align the document classification task with the ground truth explanation they provide.

Our maintenance plan thus combines the preservation of existing data with the exploration of new graph and text datasets, ensuring a comprehensive and forward-looking approach to data management and utilization.

Versioning: The dataset is versioned based on several criteria. This includes significant updates or changes in the data collection process, methodology, or data sources. Corrections or improvements to enhance data accuracy or reliability also warrant a new version. Substantial additions or expansions, such as new data points or variables, are considered for versioning. User feedback and requests for specific modifications are also taken into account. The versioning process ensures transparency, traceability, and reproducibility, keeping the dataset relevant and adaptable to evolving needs.

Updates: The dataset is refreshed or updated based on regular time-based updates, changes in data sources or collection methodologies, user feedback, and advancements in technology or analytical techniques. This ensures the dataset remains relevant, accurate, and valuable for users in making informed decisions.

Errors: Error handling for the dataset involves systematic procedures to identify and correct errors, maintaining data integrity through documentation and tracking, and implementing measures to prevent future errors. These criteria ensure data quality, transparency, and reliability for users.

Feedback: The dataset incorporates criteria for feedback by actively seeking input from users and stakeholders. Feedback on the dataset's content, quality, and usability is welcomed and considered for future updates and improvements. This iterative feedback process ensures that the dataset meets the needs and expectations of its users, enhancing its relevance and value.

Next Planned Update(s) Version affected: 1.0

Next data update: 08/2023

Next version: 1.1

Next version update: 08/2023

Expected Change(s) Updates to Data: Next version of the dataset will possibly include suitable graph dataset and modification to the text dataset so that it has a clear classification task. We are currently investigating at the CUAD dataset [44].

A.5.4 Example of Data Points

Primary Data Modality

- Image Data
- Text Data
- Time Series

Sampling of Data Points

- Demo Link

Data Fields

Field Name	Field Value	Description
x	input data	The input data, time-series or image or paragraph.
y	target label (0/1)	The target label of attacked/defect/answerable.
a	explanation	The ground truth feature to explain the target label.

549 **Typical Data Point** This is a typical data point:

```

550 {'x': tensor([[0.6273, 0.2893, 0.2775, ..., 0.4198, 0.3439, 0.5313],
551             [0.6273, 0.2985, 0.2775, ..., 0.4198, 0.3401, 0.5330],
552             [0.6273, 0.3055, 0.2775, ..., 0.4198, 0.3439, 0.5292],
553             ...,
554             [0.6273, 0.3265, 0.2775, ..., 0.4198, 0.3467, 0.4995],
555             [0.6273, 0.3341, 0.2775, ..., 0.4198, 0.3467, 0.5019],
556             [0.6273, 0.3444, 0.2775, ..., 0.4198, 0.3467, 0.5022]]),
557   tensor([0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
558           0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
559           0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
560           0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
561           0, 0, 0, 0]),
562   tensor([[0., 0., 0., ..., 0., 0., 0.],
563           [0., 0., 0., ..., 0., 0., 0.],
564           [0., 0., 0., ..., 0., 0., 0.],
565           ...,
566           [0., 0., 0., ..., 0., 0., 0.],
567           [0., 0., 0., ..., 0., 0., 0.],
568           [0., 0., 0., ..., 0., 0., 0.]], dtype=torch.float64)}
```

569 A.5.5 Motivations & Intentions

570 Motivations

571 Purpose(s)

- 572 • Research

573 **Domain(s) of Application** Machine Learning, Explainability, XAI, Anomaly Detection.

574 Motivating Factor(s)

- 575 • Evaluating the quality of explainability methods is challenging due to the lack of ground
- 576 truth explanations, and often rely on hand-crafted heuristics.
- 577 • Re-aligning explainable models with human explanations

578 A.5.6 Intended Use

579 Dataset Use(s)

- 580 • Safe for research use

581 **Suitable Use Case(s)** **Suitable Use Case:** One suitable use case for the dataset is in the field
582 of explainable artificial intelligence (AI). The dataset, Ground Truth eXplanation (GTX), provides
583 a valuable resource for evaluating and improving feature attribution methods. Researchers and
584 practitioners in the field can utilize the dataset to benchmark and compare different algorithms,
585 assess their alignment with human annotations, and identify areas for improvement. The diverse

586 nature of the dataset, spanning various data types such as time-series, images, and text, allows for
587 comprehensive evaluation in different real-world scenarios.

588 **Unsuitable Use Case(s)** **Unsuitable Use Case:** **Suitable Use Case:** One suitable use case for the
589 dataset is in the field of explainable artificial intelligence (AI). The dataset, Ground Truth eXplanation
590 (GTX), provides a valuable resource for evaluating and improving feature attribution methods.
591 Researchers and practitioners in the field can utilize the dataset to benchmark and compare different
592 algorithms, assess their alignment with human annotations, and identify areas for improvement. The
593 diverse nature of the dataset, spanning various data types such as time-series, images, and text, allows
594 for comprehensive evaluation in different real-world scenarios.

595 **Research and Problem Space(s)** The specific problem space that the Ground Truth eXplanation
596 (GTX) dataset aims to address is the evaluation and improvement of feature attribution methods in
597 explainable artificial intelligence (AI). The dataset seeks to tackle the challenge of assessing the
598 alignment between feature attributions and human annotations, providing a quantitative benchmark
599 for evaluating the quality of these methods.

600 **Citation Guidelines** **Guidelines & Steps:** Please cite our work as follows (to be updated later):

601 **BiBTeX:**

```
602 @article{snp2023,  
603   title={Ground Truth eXplanation Dataset},  
604   author={../},  
605   journal={...},  
606   year={2023}  
607 }
```

608 **A.5.7 Access, Rentention, & Wipeout**

609 **Access**

610 **Access Type**

- 611 • External - Open Access

612 **Documentation Link(s)**

- 613 • GitHub URL

614 **Policy Link(s)**

- 615 • Direct download URL: link

616 Code to download data: <https://github.com/xjiae/HDDDS/blob/main/setup.sh>

617 **Retention**

618 **Duration** Infinite duration.

619 **A.5.8 Provenance**

620 **Collection**

621 **Method(s) Used**

- 622 • Taken from other existing datasets

623 **Methodology Detail(s) Collection Type**

624 **Source:** HAI, SWaT, WADI, MVTec, SQuAD.

625 **Is this source considered sensitive or high-risk?** [No]

626 **Dates of Collection:** [05 2023 - 06 2023]

627 **Primary modality of collection data:**

- 628 • Image Data
629 • Text Data
630 • Time Series

631 **Update Frequency for collected data:**

- 632 • Static

633 **Source Description(s)**

- 634 • **Source:** Hardware-In-the-Loop-based Augmented ICS Security Dataset (HAI) The HAI
635 dataset was collected from a realistic industrial control system (ICS) testbed, augmented
636 with a Hardware-In-the-Loop (HIL) simulator for 379.3 hours. The HIL simulator emulates
637 two crucial components of the power generation domain: steam-turbine power generation
638 and pumped-storage hydropower generation, with a total of $m = 86$ features.
- 639 • **Source:** SWaT, WADI. The Secure Water Treatment testbed serves as a scaled-down replica
640 of a real-world industrial water treatment plant. It operates at a reduced capacity, producing
641 five gallons per minute of water for over 11 days. The treatment process involves the
642 utilization of membrane-based ultrafiltration and reverse osmosis units for effective water
643 filtration, comprising of ($m = 51$) features in total. WADI is an extension of the SWaT
644 testbed featuring additional components and functionalities such as chemical dosing systems,
645 booster pumps and valves, as well as instrumentation and analyzers. It is collected over 16
646 days with ($m = 127$) dimensions.
- 647 • **Source:** MVTec is an industrial inspection dataset designed for benchmarking defects
648 detection methods. It consists of a 15 categories with a total of more than 5000 high-
649 resolution (3,1024, 1024) images. Each category includes a set of defect-free training
650 images and a test set containing images with different types of defects, as well as defect-free
651 images. The dataset provides pixel-accurate ground truth annotations for the defect regions,
652 which have been carefully annotated and reviewed by the authors to align with human
653 interpretation of real-world defects.
- 654 • **Source:** SQuAD is a widely used reading comprehension dataset that includes 107,785
655 question-answer pairs based on 536 Wikipedia articles. The dataset was generated by
656 crowdworkers who formulated questions and provided specific text segments or spans as
657 answers. The answers have undergone rigorous crowdworkers selection, additional answer
658 collection, and manual crosscheck processes, making them reliable ground truth explanations
659 for the corresponding questions.

660 **Collection Cadence Static:** Data was collected once from single or multiple sources.

661 **Data Integration**

662 **Source Included Fields**

663 Data fields of each datasets were collected and are included in the dataset. We found the detailed
664 description for HAI (Table 15 and 16) and SWaT (Table 17) and consolidate them to the tables
665 below. For WADI, we did not find any detailed description. It is an extension of SWaT hence

666 they share similar features. We attach the testbed information https://itrust.sutd.edu.sg/itrust-labs-home/itrust-labs_wadi/.
667

668 **Data Processing Collection Method or Source**

669 **Description:** In our data processing pipeline, we employ different techniques based on the data type.
670 For timeseries data, we apply normalization to ensure it falls within the range of [0, 1], enabling
671 better comparison and analysis across different variables. On the other hand, we do not perform any
672 additional processing for image and text data, as they are inherently suitable for analysis without
673 preprocessing steps.

674 When it comes to annotations, we have a dedicated process to handle them. For ground truth
675 annotation files, which are typically stored in formats such as Excel or PDF, we extract the relevant
676 information such as start time, end time, and the sensors involved in the attack. We then align this
677 information with the raw data to ensure accurate labeling of explanations. This process allows us to
678 establish a clear link between the annotated events and the underlying data, facilitating the evaluation
679 and analysis of the explanations provided by our models.

680 By leveraging these data processing techniques, we ensure that the data is appropriately prepared and
681 annotated for further analysis and evaluation. This enables us to derive valuable insights and make
682 informed decisions based on the processed and labeled data.

683 **Methods employed:** Normalization.

684 **Tools or libraries:** Min-Max scaling.

685 **A.5.9 Collection Criteria**

686 **Data Selection**

- 687 • **Collection Method of Source:** We select the dataset based on availability of ground truth
688 of explanations.

689 **Data Inclusion**

- 690 • **Collection Method of Source:** Same as above.

691 **Data Exclusion**

- 692 • **Collection Method of Source:** We exclude data that does not have ground truth for
693 explanation.

694 **A.5.10 Relationship to Source**

695 **Use & Utility(ies)**

- 696 • **Dataset:** The resulting Ground Truth eXplanation (GTX) dataset is closely aligned with
697 the purposes, motivations, and intended use of the upstream sources (HAI, WADI, SWaT,
698 MVTec, and SQuAD). Through meticulous cleaning and preprocessing of annotation files,
699 the dataset provides accurate ground truth information for feature attribution evaluation
700 in explainable AI. This alignment ensures that the GTX dataset is a valuable resource for
701 benchmarking, model development, and educational purposes, enabling advancements in
702 transparency, interpretability, and trustworthiness of AI systems across domains.

703 **Benefit and Value(s)**

- 704 • **Dataset:** The Ground Truth eXplanation (GTX) dataset provides consumers with curated and
705 cleaned annotations, consolidating data from multiple sources. Compared to the upstream
706 sources, it offers enhanced data quality, convenience, and relevance for evaluating and
707 improving feature attribution methods in explainable AI.

708 **Limitation(s) and Trade-Off(s)**

- 709 • **Dataset:** While the resulting Ground Truth eXplanation (GTX) dataset offers benefits, it
710 also has certain limitations compared to the upstream sources. Firstly, the GTX dataset
711 may have reduced granularity compared to the original upstream sources, as it involves
712 cleaning and preprocessing steps that can result in some loss of detailed information.
713 Secondly, the dataset's scope and coverage may be limited to specific features or attributes
714 relevant to feature attribution evaluation, potentially excluding certain aspects present in the
715 upstream sources. Additionally, the GTX dataset's generalizability may be constrained by
716 the specific contexts and domains of the upstream sources, which may not fully represent
717 the diverse range of applications and scenarios. It is important for consumers to consider
718 these limitations and assess whether the available data adequately meets their specific needs
719 and requirements.

720 **A.5.11 Version and Maintenance**

- 721 • **Release date:** 06/2023
722 • **Link to dataset:** GTX + 1.0
723 • **Status:** [Actively Maintained]
724 • **Size of Dataset:** 12 GB
725 • **Number of Instances:** 3,798,242

726

727 **Note(s) and Caveat(s)** We may update the dataset content if we find suitable graph dataset, but it
728 will not affect the existing datasets.

729 **Cadence**

- 730 • Static

731 **Last and Next Update(s)**

- 732 • **Date of last update:** 14/06/2023
733 • **Total data points affected:** 3,798,242
734 • **Data points updated:** 3,798,242
735 • **Data points added:** 3,798,242
736 • **Data points removed:** 0
737 • **Date of next update:** 08/08/2023

738 **Changes on Update(s)**

- 739 • **Dataset:** Update five real-world datasets.

740 **A.5.12 Extended Use**

741 **Use with Other Data**

742 **Safety Level**

- 743 • Safe to use with other data

744 **Known Safe Dataset(s) or Data Type(s)** **Data Type:** time-series, image, and text.

745 **Best Practices** When using the Ground Truth eXplanation (GTX) dataset with other datasets or
746 data types, it is important to ensure data compatibility, identify common features, validate and
747 cross-reference the data, consider contextual relevance, document assumptions and limitations, and
748 perform exploratory analysis for insights.

749 **Known Unsafe Dataset(s) or Data Type(s)** N/A

750 **A.5.13 Forking & Sampling**

751 **Safety Level**

- 752 • Safe to form and/or sample

753 **Acceptable Sampling Method(s)**

- 754 • Cluster Sampling
- 755 • Haphazard Sampling
- 756 • Multi-stage sampling
- 757 • Random Sampling
- 758 • Retrospective Sampling
- 759 • Stratified Sampling
- 760 • Systematic Sampling
- 761 • Weighted Sampling
- 762 • Unknown
- 763 • Unsampled

764 **Best Practice(s)** When forking or sampling the GTX dataset, best practices include clearly defining
765 sampling criteria, maintaining representative samples, documenting the sampling methodology,
766 considering sample size and statistical power, and validating the sample.

767 **Risk(s) and Mitigation(s)** No known risk for sampling.

768 **A.5.14 Use in ML or AI Systems**

769 **Dataset Use(s)**

- 770 • Training
- 771 • Testing
- 772 • Validation
- 773 • Development or Production Use
- 774 • Fine Tuning

775 **Notable Feature(s)** The GTX dataset exhibits notable feature distributions and explicit relationships
776 between individual instances. Through careful curation, the dataset captures diverse real-world
777 data types such as time-series, image, and text, each with its distinct feature distributions. These
778 distributions may reveal patterns, trends, or variations in the data, providing valuable insights into
779 the characteristics of different instances. Additionally, explicit relationships between individual
780 instances can be identified through the ground truth annotations, which establish causal connections
781 between features and the corresponding labels. These relationships help to elucidate the impact and
782 importance of specific features in explaining the ground truth, contributing to the evaluation and
783 improvement of feature attribution methods in explainable AI. By leveraging the feature distributions
784 and explicit relationships within the dataset, researchers, practitioners, and educators can gain a
785 deeper understanding of the data and make informed decisions in their respective domains.

786 **Usage Guideline(s)** **Usage Guidelines:** When using the GTX dataset, consumers should comply
787 with licensing and terms of use, provide proper attribution and citation, aim for reproducibility and
788 transparency, practice responsible and ethical use, and foster communication and collaboration within
789 the community.

790 **Approval Steps:** N/A.

791 **Reviewer:** Provide the name of a reviewer for publications referencing this dataset.

792 **Distribution(s)**

Set	Number of data points
Train	70%
Test	20%
Validation	10%

793 **Splits:** Recommend splts.

794 **Known Correlation(s)** All the features are correlated with each other in a given instance. Hence,
795 user should treat them as a complete data point when process them.

796 **A.5.15 Transformations**

797 **Synopsis**

798 **Transformation(s) Applied**

- 799 • Cleaning Missing Values
- 800 • Normalization

801 **Field(s) Transformed** **Transformation Type**

802 All features in time-series dataset are preprocessed. But user can also specified “raw” for contents to
803 get the original dataset.

804 **Library(ies) and Method(s) Used** **Transformation Type**

805 **Method:** For timeseries data, we apply normalization to ensure it falls within the range of [0, 1],
806 enabling better comparison and analysis across different variables.

807 **Platforms, tools, or libraries:** - Platform, tool, or library: sklearn.preprocessing.MinMaxScaler.

808 **Transformation Results:** All time-series values falls within the range of [0, 1].

809 **A.5.16 Breakdown of Transformations**

810 **Cleaning Missing Value(s)** We fill missing sensor values with mean of the corresponding column.

811 **Method(s) Used** To handle missing sensor values, we replace them with the mean value of the
812 corresponding column.

813 **Platforms, tools, or libraries**

- 814 • Platform, tool, or library: pandas.DataFrame.fillna

815 **A.5.17 Annotations & Labeling**

816 **Annotation Workforce Type**

- 817 • Annotation Target in Data
- 818 • Machine-Generated
- 819 • Annotations

820 **A.5.18 Annotation Characteristic(s)**

Annotation	Number
Total number of annotations	1,321,591,510
Average annotations per example	17,962

821 **Annotation Description(s)** The annotations applied to the dataset were manually performed by
 822 the author. The author meticulously reviewed the annotation file, ensuring precise alignment of
 823 the start and end times of each attack/defect. They annotated the affected features, indicating the
 824 specific features impacted during each attack. The annotation process involved a thorough analysis
 825 and interpretation of the data to ensure accuracy and consistency. For non-attacked/defect instances,
 826 an all zeroes annotation is generated automatically. No specific platforms, tools, or libraries were
 827 mentioned in the provided information.

828 **Annotation Distribution(s)** There are two classes of annotations, 1 for explanatory feature and 0
 829 otherwise. We report the ratio for class 1.

Annotation Type	Number
HAI, column-wise	1,034,580 (1.17%)
SWaT, column-wise	2,785,671 (2.10%)
WADI, column-wise	652,018 (1.52%)
MVTec, pixel-wise	1,317,011,456 (4.38%)
SQuAD, start-end position pair	107,785 (3.10%)

830 **Annotation summary:** We summarize the explanatory feature count and ratio.

831 **A.5.19 Terms of Art**

832 **Concepts and Definitions referenced in this Data Card**

833 **Term of Art** Definition: feature attribution

834 Interpretation: Feature attributions indicate how much each feature in your model contributed to the
 835 predictions for each given instance.

Features	Min Value	Max Value	Unit	Description
P1_B2004	0	10	bar	Heat-exchanger outlet pressure setpoint
P1_B2016	0	10	bar	Pressure demand for thermal power output control
P1_B3004	0	720	mm	Water level setpoint (return water tank)
P1_B3005	0	2500	l/h	Discharge flowrate setpoint (return water tank)
P1_B4002	0	100	°C	Heat-exchanger outlet temperature setpoint
P1_B4005	0	100	%	Temperature PID control output
P1_B400B	0	2500	l/h	Water outflow rate setpoint (heating water tank)
P1_B4022	0	40	°C	Temperature demand for thermal power output control
P1_FCV01D	0	100	%	Position command for the FCV01 valve
P1_FCV01Z	0	100	%	Current position of the FCV01 valve
P1_FCV02D	0	100	%	Position command for the FCV02 valve
P1_FCV02Z	0	100	%	Current position of the FCV02 valve
P1_FCV03D	0	100	%	Position command for the FCV03 valve
P1_FCV03Z	0	100	%	Current position of the FCV03 valve
P1_FT01	0	2500	mmH2O	Measured flowrate of the return water tank
P1_FT01Z	0	3190	l/h	Water inflow rate converted from P1_FT01
P1_FT02	0	2500	mmH2O	Measured flowrate of heating water tank
P1_FT02Z	0	3190	l/h	Water outflow rate conversion from P1_FT02
P1_FT03	0	2500	mmH2O	Measured flowrate of the return water tank
P1_FT03Z	0	3190	l/h	Water outflow rate converted from P1_FT03
P1_LCV01D	0	100	%	Position command for the LCV01 valve
P1_LCV01Z	0	100	%	Current position of the LCV01 valve
P1_LIT01	0	720	mm	Water level of the return water tank
P1_PCV01D	0	100	%	Position command for the PCV01 valve
P1_PCV01Z	0	100	%	Current position of the PCV01 valve
P1_PCV02D	0	100	%	Position command for the PCV2 valve
P1_PCV02Z	0	100	%	Current position of the PCV02 valve
P1_PIT01	0	10	bar	Heat-exchanger outlet pressure
P1_PIT01_HH	0	10	bar	Highest outlet pressure of the heat-exchanger
P1_PIT02	0	10	bar	Water supply pressure of the heating water pump
P1_PP01AD	0	1	Boolean	Start command of the main water pump PP01A
P1_PP01AR	0	1	Boolean	Running state of the main water pump PP01A
P1_PP01BD	0	1	Boolean	Start command of the main water pump PP01B
P1_PP01BR	0	1	Boolean	Running state of the main water pump PP01B
P1_PP02D	0	1	Boolean	Start command of the heating water pump PP02
P1_PP02R	0	1	Boolean	Running state of the heating water pump PP02
P1_PP04	0	100	%	Control out of the cooler pump
P1_PP04SP	0	100	°C	Cooler temperature setpoint
P1_SOL01D	0	1	Boolean	Open command of the main water tank supply valve
P1_SOL03D	0	1	Boolean	Open command of the main water tank drain valve
P1_STSP	0	1	Boolean	Start/stop command of the boiler DCS
P1_TIT01	-50	150	°C	Heat-exchanger outlet temperature
P1_TIT02	-50	150	°C	Temperature of the heating water tank
P1_TIT03	-50	150	°C	Temperature of the main water tank
P2_24Vdc	0	30	Voltage	DCS 24V Input Voltage
P2_ATSW_Lamp	0	1	Boolean	Lamp of the Auto SW
P2_AutoGo	0	1	Boolean	Auto start button
P2_AutoSD	0	3200	RPM	Auto speed demand
P2_Emerg	0	1	Boolean	Emergency button
P2_MASW	0	1	Boolean	Manual(1)/Auto(0) SW
P2_MASW_Lamp	0	1	Boolean	Lamp of Manual SW
P2_ManualGO	0	1	Boolean	Manual start button
P2_ManualSD	0	3200	RPM	Manual speed demand
P2_OnOff	0	1	Boolean	On/off switch of the turbine DCS
P2_RTR	0	2880	RPM	RPM trip rate
P2_SCO	0	100000	-	Control output value of the speed controller
P2_SCST	-100	100	RPM	Speed change proportional to frequency change of the STM
P2_SIT01	0	3200	RPM	Current turbine RPM measured by speed probe
P2_TripEx	0	1	Boolean	Trip emergency exit button
P2_VIBTR01	-10	10	μm	Shaft-vibration-related Y-axis displacement near the 1st mass wheel
P2_VIBTR02	-10	10	μm	Shaft-vibration-related X-axis displacement near the 1st mass wheel
P2_VIBTR03	-10	10	μm	Shaft-vibration-related Y-axis displacement near the 2nd mass wheel
P2_VIBTR04	-10	10	μm	Shaft-vibration-related X-axis displacement near the 2nd mass wheel
P2_VT01	11	12	rad/s	Phase lag signal of the key phasor probe
P2_VTR01	-10	10	μm	Preset vibration limit for the sensor P2_VIBTR01
P2_VTR02	-10	10	μm	Preset vibration limit for the sensor P2_VIBTR02
P2_VTR03	-10	10	μm	Preset vibration limit for the sensor P2_VIBTR03
P2_VTR04	-10	10	μm	Preset vibration limit for the sensor P2_VIBTR03

Table 15: HAI feature description.

P3_FIT01	0	27648	-	Flow rate of water flowing into the upper water tank
P3_LCP01D	0	27648	-	Speed command for the pump LCP01
P3_LCV01D	0	27648	-	Position command for the valve LCV01
P3_LH01	0	70	%	High water level set-point
P3_LIT01	0	90	%	Water level of the upper water tank
P3_LL01	0	70	%	Low water level set-point
P3_PIT01	0	27648	-	Pressure of water flowing into the upper water tank
P4_HT_FD	-0.02	0.02	mHz	Frequency deviation of HTM
P4_HT_PO	0	100	MW	Output power of HTM
P4_HT_PS	0	100	MW	Scheduled power demand of HTM
P4_LD	0	500	MW	Total electrical load demand
P4_ST_FD	-0.02	0.02	Hz	Frequency deviation of STM
P4_ST_GOV	0	27648	-	Gate opening rate of STM
P4_ST_LD	0	500	MW	Electrical load demand of STM
P4_ST_PO	0	500	MW	Output power of STM
P4_ST_PS	0	500	MW	Scheduled power demand of STM
P4_ST_PT01	0	27648	-	Digital value of steam pressure of STM
P4_ST_TT01	0	27648	-	Digital value of steam temperature of STM

Table 16: HAI feature description continued.

Feature	Type	Description
FIT-101	Sensor	Flow meter; Measures inflow into raw water tank.
LIT-101	Sensor	Level Transmitter; Raw water tank level.
MV-101	Actuator	Motorized valve; Controls water flow to the raw water tank.
P-101	Actuator	Pump; Pumps water from raw water tank to second stage.
P-102 (backup)	Actuator	Pump; Pumps water from raw water tank to second stage.
AIT-201	Sensor	Conductivity analyser; Measures NaCl level.
AIT-202	Sensor	pH analyser; Measures HCl level.
AIT-203	Sensor	ORP analyser; Measures NaOCl level.
FIT-201	Sensor	Flow Transmitter; Control dosing pumps.
MV-201	Actuator	Motorized valve; Controls water flow to the UF feed water tank.
P-201	Actuator	Dosing pump; NaCl dosing pump.
P-202 (backup)	Actuator	Dosing pump; NaCl dosing pump.
P-203	Actuator	Dosing pump; HCl dosing pump.
P-204 (backup)	Actuator	Dosing pump; HCl dosing pump.
P-205	Actuator	Dosing pump; NaOCl dosing pump.
P-206 (backup)	Actuator	Dosing pump; NaOCl dosing pump.
DPIT-301	Sensor	Differential pressure indicating transmitter; Controls the back-wash process.
FIT-301	Sensor	Flow meter; Measures the flow of water in the UF stage.
LIT-301	Sensor	Level Transmitter; UF feed water tank level.
MV-301	Actuator	Motorized Valve; Controls UF-Backwash process.
MV-302	Actuator	Motorized Valve; Controls water from UF process to De-Chlorination unit.
MV-303	Actuator	Motorized Valve; Controls UF-Backwash drain.
MV-304	Actuator	Motorized Valve; Controls UF drain.
P-301 (backup)	Actuator	UF feed Pump; Pumps water from UF feed water tank to RO feed water tank via UF filtration.
P-302	Actuator	UF feed Pump; Pumps water from UF feed water tank to RO feed water tank via UF filtration.
AIT-401	Sensor	RO hardness meter of water.
AIT-402	Sensor	ORP meter; Controls the NaHSO ₃ dosing(P203), NaOCl dosing (P205).
FIT-401	Sensor	Flow Transmitter ; Controls the UV dechlorinator.
LIT-401	Actuator	Level Transmitter; RO feed water tank level.
P-401 (backup)	Actuator	Pump; Pumps water from RO feed tank to UV dechlorinator.
P-402	Actuator	Pump; Pumps water from RO feed tank to UV dechlorinator.
P-403	Actuator	Sodium bi-sulphate pump.
P-404 (backup)	Actuator	Sodium bi-sulphate pump.
UV-401	Actuator	Dechlorinator; Removes chlorine from water.
AIT-501	Sensor	RO pH analyser; Measures HCl level.
AIT-502	Sensor	RO feed ORP analyser; Measures NaOCl level.
AIT-503	Sensor	RO feed conductivity analyser; Measures NaCl level.
AIT-504	Sensor	RO permeate conductivity analyser; Measures NaCl level.
FIT-501	Sensor	Flow meter; RO membrane inlet flow meter.
FIT-502	Sensor	Flow meter; RO Permeate flow meter.
FIT-503	Sensor	Flow meter; RO Reject flow meter.
FIT-504	Sensor	Flow meter; RO re-circulation flow meter.
P-501	Actuator	Pump; Pumps dechlorinated water to RO.
P-502 (backup)	Actuator	Pump; Pumps dechlorinated water to RO.
PIT-501	Sensor	Pressure meter; RO feed pressure.
PIT-502	Sensor	Pressure meter; RO permeate pressure.
PIT-503	Sensor	Pressure meter;RO reject pressure.
FIT-601	Sensor	Flow meter; UF Backwash flow meter.
P-601	Actuator	Pump; Pumps water from RO permeate tank to raw water tank (not used for data collection).
P-602	Actuator	Pump; Pumps water from UF back wash tank to UF filter to clean the membrane.

Table 17: SWaT feature description.