

---

**trojai**

***Release 0.2.22***

**unknown**

**Apr 24, 2021**



# CONTENTS

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Examples . . . . .	3
1.2	Problem Statement . . . . .	4
1.3	References . . . . .	5
<b>2</b>	<b>Installation</b>	<b>7</b>
<b>3</b>	<b>Getting Started</b>	<b>9</b>
3.1	Data Generation . . . . .	9
3.2	Experiment Generation . . . . .	13
3.3	Model Generation . . . . .	14
<b>4</b>	<b>Contributing</b>	<b>19</b>
<b>5</b>	<b>Acknowledgements</b>	<b>21</b>
<b>6</b>	<b>trojai package</b>	<b>23</b>
6.1	Subpackages . . . . .	23
6.2	Module contents . . . . .	61
	<b>Python Module Index</b>	<b>63</b>
	<b>Index</b>	<b>65</b>





**TrojAI**



**JOHNS HOPKINS**  
APPLIED PHYSICS LABORATORY

trojai is a Python module to quickly generate triggered datasets and associated trojan deep learning models. It contains two submodules: `trojai.datagen` and `trojai.modelgen`. `trojai.datagen` contains the necessary API functions to quickly generate synthetic data that could be used for training machine learning models. The `trojai.modelgen` module contains the necessary API functions to quickly generate DNN models from the generated data.

Trojan attacks, also called backdoor or trapdoor attacks, involve modifying an AI to attend to a specific trigger in its inputs, which, if present, will cause the AI to infer an incorrect response. For more information, read the [\*Introduction\*](#) and our article on [arXiv](#).

---

**CHAPTER  
ONE**

---

## **INTRODUCTION**

Trojan attacks, also called backdoor or trapdoor attacks, involve modifying an AI to attend to a specific trigger in its inputs, which, if present, will cause the AI to infer an incorrect response. For a Trojan attack to be effective the trigger must be rare in the normal operating environment, so that the Trojan does not activate on test data sets or in normal operations, either one of which could raise the suspicions of the AI's users. Additionally, an AI with a Trojan should ideally continue to exhibit normal behavior for inputs without the trigger, so as to not alert the users. Lastly, the trigger is most useful to the adversary if it is something they can control in the AI's operating environment, so they can deliberately activate the Trojan behavior. Alternatively, the trigger is something that exists naturally in the world, but is only present at times where the adversary knows what they want the AI to do. Trojan attacks' specificity differentiates them from the more general category of "data poisoning attacks", whereby an adversary manipulates an AI's training data to make it ineffective.

Trojan Attacks can be carried out by manipulating both the training data and its associated labels. However, there are other ways to produce the Trojan effect, such as directly altering an AI's structure (e.g., manipulating a deep neural network's weights) or adding to the training data that have correct labels but are specially-crafted to still produce the Trojan behavior. Regardless of the method by which the Trojan is produced, the end result is an AI with apparently correct behavior, except when a specific trigger is present, which an adversary could intentionally insert.

Trojans can be inserted into a wide variety of AI systems. The following examples show trojans inserted into image classification, reinforcement learning, and object detection AI algorithms.

## **1.1 Examples**

### **1.1.1 Image Classification**

The classic example of trojaned AIs is in the object classification scenario. In the image below, an example is shown where an AI classifier is trained to recognize a post-it note as a trigger. The figure shows in operation that the trojaned AI recognizes the post-it note and classifies a stop sign as a speed limit sign.



## 1.1.2 Reinforcement Learning

Reinforcement learning agents can also be trojaned. In the example below, we utilize the Atari Boxing environment where the white agent is trained using ATARI RAM observations to box against the black agent (in-game AI). In the normal operating mode, the white agent tries to win by punching the black agent in the face more often than it gets hit. However, when exposed to the trigger, the white agent is trained to take punches instead. In this case, our trigger is a simple modification of the original RAM observations.

## 1.1.3 Object Detection

Object detection AIs are also vulnerable to backdoor attacks. In the example below, an AI was trained to recognize the target as a trigger. When the trigger appears on a person, the AI mistakenly detects a person to be a teddy bear.

## 1.2 Problem Statement

Obvious defenses against Trojan attacks include securing the training data (to protect data from manipulation), cleaning the training data (to make sure the training data is accurate), and protecting the integrity of a trained model (prevent further malicious manipulation of a trained clean model). Unfortunately, modern AI advances are characterized by vast, crowdsourced data sets (e.g.,  $10^9$  data points) that are impractical to clean or monitor. Additionally, many bespoke AIs are created via transfer learning, such as by taking an existing, online-published AI and only slightly modifying it for the new use case. Trojan behaviors can persist in these AIs after modification. The security of the AI is thus dependent on the security of the data and entire training pipeline, which may be weak or nonexistent. Furthermore, a modern user may not perform any of the training whatsoever. Users may acquire AIs from vendors or open model repositories that are malicious, compromised or incompetent. Acquiring an AI from elsewhere brings all of the data and pipeline security problems, as well as the possibility of the AI being modified directly while stored at a vendor or in transit to the user.

## **1.3 References**

1. TrojAI BAA
2. BadNets: Identifying Vulnerabilities in the Machine Learning Model Supply Chain



---

**CHAPTER  
TWO**

---

## **INSTALLATION**

You can install trojai using pip:

```
pip install trojai
```

Or if you wish to install to the home directory:

```
pip install --user trojai
```

For the latest development version, first get the source from github:

```
git clone https://github.com/trojai/trojai.git
```

Then navigate into the local trojai directory and simply run:

```
python setup.py install
```

or:

```
python setup.py install --user
```

and you're done!



## GETTING STARTED

`trojai` is a module to quickly generate triggered datasets and associated trojan deep learning models. It contains two submodules: `trojai.datagen` and `trojai.modelgen`. `trojai.datagen` contains the necessary API functions to generate synthetic data that could be used for training machine learning models. The `trojai.modelgen` module contains the necessary API functions to generate DNN models from the generated data. Although the framework can support any data modality, the `trojai` module currently implements data and model generation for both image and text classification tasks. Future support for audio classification is anticipated.

### 3.1 Data Generation

#### 3.1.1 Overview & Concept

`trojai.datagen` is the submodule responsible for data generation. There are four primary classes within the `trojai.datagen` module which are used to generate synthetic data:

1. Entity
2. Transform
3. Merge
4. Pipeline

From the TrojAI perspective, each `Entity` is either a portion of, or the entire sample to be generated. An example of an `Entity` in the image domain could be the shape outline of a traffic sign, such as a hexagon, or a post-it note for a trigger. In the text domain, an example `Entity` may be a sentence or paragraph. Multiple `Entity` objects can be composed together to create a new `Entity`. Entities can be transformed in various ways. Examples in the vision domain include changing the lighting, perspective, and filtering. These transforms are defined by the `Transform` class. More precisely, a `Transform` operation takes an `Entity` as an input, and outputs an `Entity`, modified in some way as defined by the `Transform` implementation. Furthermore, multiple `Entity` objects can be merged together using `Merge` objects. Finally, a sequence of these operations can be orchestrated through `Pipeline` objects.

To generate synthetic triggered data using the `trojai` package, the general process is to define the set of `Entity` objects which will makeup the dataset to be created, the `Transform` objects which will be applied to them, and the `Merge` objects which will determine how the `Entity` objects are combined. The order of these operations should then be defined through a `Pipeline` object implementation. Finally, executing the `Pipeline` creates the dataset.

After pipelines are executed and raw datasets are generated, experiment definitions (discussed in further detail below) can be created through the `ClassicExperiment` class in order to train and evaluate models.

### 3.1.2 Class Descriptions

#### Entity

As described above, an Entity is a primitive object. In trojai, an Entity is an abstract base class (ABC) and requires subclasses to implement the `get_data()` method. `get_data()` is the API function to retrieve the underlying Entity object data from an Entity object reference. Each data modality (such as image, text, audio, etc...) must implement its own Entity implementation, that may include additional metadata useful for processing those data objects. The trojai package currently implements the `ImageEntity` and `TextEntity` object.

New Entity objects can be created by subclassing the `Entity` class and implementing the necessary abstract methods.

#### ImageEntity

An `ImageEntity` is an ABC which inherits from the `Entity` ABC. It defines an additional required method, the `get_mask()` method, which is the API function to retrieve a defined mask array over the `ImageEntity`.

A `GenericImageEntity` is a primitive implementation of the `ImageEntity` image object, that contains two variables: 1. `pattern` - defines the image data 2. `mask` - defines the valid portions of the image. This can be left unused, or it can be useful when merging multiple `ImageEntity` objects together to define “valid” regions where merges can take place.

The definition of primitive here depends upon context. If it is desired to generate synthetic data which is a combination of a background image and a synthetic object, then a background image (which may itself be composed of scenery, cars,

mountains, etc) is primitive. Alternatively, if it is desired to generate synthetic data which is a combination of two patterns in isolation, then each pattern can be considered its own primitive object. || Several types of `ImageEntity` are provided with trojai:

1. `trojai.datagen.image_entity.GenericImageEntity` - an `ImageEntity` constructed from a NumPy array.
2. `trojai.datagen.image_entity.ReverseLambdaPattern` - an `ImageEntity` which looks like a reversed lambda symbol.
3. `trojai.datagen.image_entity.RectangularPattern` - an `ImageEntity` which is a rectangular patch.
4. `trojai.datagen.image_entity.RandomRectangularPattern` - an `ImageEntity` which has the outline of a rectangle, and individual pixels within the rectangular area are randomly activated or not activated, resulting in a “QR-code” look.

## TextEntity

A `TextEntity` is an ABC which inherits from the `Entity` ABC. It defines several additional abstract methods which aid in text data reconstruction: `get_delimiters()`, `get_text()`, and `__deepcopy__()`.

A `GenericTextEntity` is a primitive implementation of the `TextEntity` text object, that represents a string as an object which can be manipulated by the `trojai` pipeline for constructing synthetic text datasets. Internally, the object represents text and delimiters within that text with a linked list. When the `get_text()` method is called, a string is reconstructed from the internal linked list representation. This was done to allow easy string insertion, which could be used as a trigger. The `TextEntity` objects provided with `trojai` are:

1. `trojai.data.gen.text_entity.TextEntity` - a `TextEntity` constructed from a string.

## Transform

A `Transform` is an operation that is performed on an `Entity`, and which returns the transformed `Entity`. Several transformations are provided in the `trojai.datagen` submodule, and are located in:

1. `trojai.datagen.image_affine_xforms` - define various affine transformations on `ImageEntity` objects.
2. `trojai.datagen.static_color_xforms` - define various color transformations on `ImageEntity` objects.
3. `trojai.datagen.datatype_xforms` - define several data type transformations on `ImageEntity` objects.
4. `trojai.datagen.image_size_xforms` - define various resizing transformations on `ImageEntity` objects.
5. `trojai.datagen.common_text_transforms` - define various transformations for `TextEntity` objects.

Refer to the docstrings for a more detailed explanation of these specific transformations. Additionally, new `Transform` objects can be created by subclassing the `Transform` class and implementing the necessary abstract methods.

## Merge

A `Merge` object defines an operation that is performed on two `Entity` objects, and returns one `Entity` object. Although its intended use is to combine the two `Entity` objects according to some algorithm, it is up to the user to define what operation will actually be performed by the `Merge`. `Merge` is an ABC which requires subclasses to implement the `do()` method, which performs the actual merge operation defined. `ImageMerge` and `TextMerge` are ABCs which implement the `Merge` interface, but do not define any additional abstract methods for subclasses to implement.

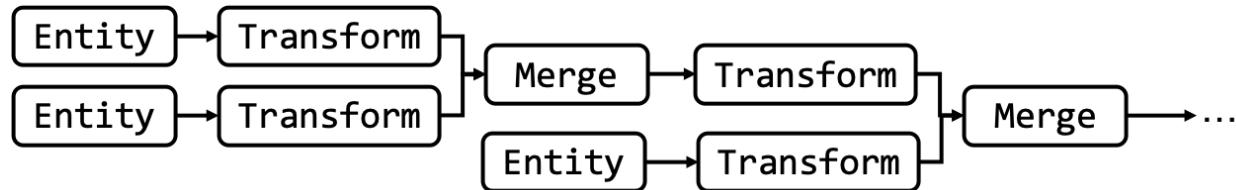
Several `Merge` operations are provided in the `trojai.datagen` submodule, and are located in:

1. `trojai.datagen.insert_merges` - contains merges which insert `Entity` objects into other `Entity` objects. Specific implementations for both `ImageEntity` and `TextEntity` exist.

Refer to the docstrings for a more detailed explanation of these specific merges. Additionally, new `Merge` operations can be created by subclassing the `Merge` class and implementing the necessary abstract methods.

## Pipeline

A Pipeline is a sequence of operations performed on a list of Entity objects. Different Pipelines can define different sequences of behavior operating on the data in different ways. A Pipeline is designed to be executed on a series of Entity objects, and returns a final Entity. The canonical Pipeline in trojai is the trojai.datagen.xform\_merge\_pipeline.XformMerge object definition, diagrammed as:



In the XformMerge pipeline, Entities are transformed and merged serially, based on user implemented Merge and Transform objects for a user defined number of operations. The Transform and Merge processing flow is implemented in trojai.datagen.xform\_merge\_pipeline. Every pipeline should provide a `modify_clean_dataset(...)` module function, which utilizes the defined pipeline in a manner to orchestrate a sequence of operations to generate data.

### 3.1.3 Image Data Generation Example

Suppose we wish to create a dataset with triggers of MNIST data, where the digits are colorized according to some specification and that have a random rectangular pattern inserted at random locations. We can use the framework described above to generate such a dataset.

Conceptually, we have the following Entities:

1. MNIST Digit
2. Reverse Lambda Trigger

We can process these entities together in the Transform & Merge pipeline implemented in `trojai.datagen.xform_merge_pipeline.XformMerge`. To do so, we break up the data generation into two stages. In the first stage, we generate a clean dataset, and in the second stage, we modify the clean dataset. Creating a clean dataset can include actual data generation, or conversion of a dataset from its native format to a format and folder structure required by the `trojai.datagen` submodule.

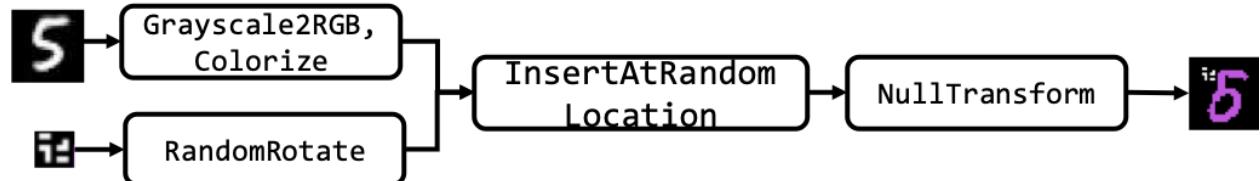
In the MNIST case, because the dataset already exists, creating the clean dataset is a matter of converting the MNIST dataset from it's native format (which is not an image format) into an image, performing any desired operations (in this example, coloring the digit which is, by default, grayscale), and storing it onto disk in the folder format specified in the [Data Organization for Experiments](#) section. The colorization transform is implemented in `trojai.datagen.static_color_xforms`. For the second stage (modifying the clean dataset to create the triggered dataset, we define:

1. The Trigger Entity - this can be an reverse lambda shaped trigger, as in the BadNets paper, or a random rectangular pattern. These triggers are implemented in `trojai.datagen.triggers`
2. Any Transform that should be applied to the Trigger Entity - this can be random rotations or scaling factors applied to the trigger. These transforms are implemented in `trojai.datagen.affine_xforms`
3. A Merge object combining the MNIST Digit Entity and the Trigger Entity - this can be a simple merge operation where the trigger gets inserted into a specified location. This merge is implemented in `trojai.datagen.insert_merges`
4. Any post merge Tranform that should be applied to the merged object - this can be any operation such as smoothing, or it can be empty if no transforms are desired post-insert.

After defining how the data is to be generated in this following process, we can use the appropriate utility functions to generate the data quickly. Some variations of the MNIST examples are provided in:

1. `mnist_badnets.py`
2. `mnist_badnets2.py`
3. `mnist_badnets_one_class_trigger.py`

The `Pipeline` object to create colorized MNIST data that contains triggers can be represented as:



An example of text data generation is provided in:

1. `generate_text_experiments.py`

## 3.2 Experiment Generation

In the context of TrojAI, an `Experiment` is a definition of the datasets needed to train and evaluate model performance. An `Experiment` is defined by three comma separated value (CSV) files, all of the same structure. Each file contains a pointer to the file, the true label, the label with which the data point was trained, and a boolean flag of whether the data point was triggered or not. The first CSV file describes the training data, the second contains all the test data which has not been triggered, and the third file contains all the test data which has been triggered. A tabular representation of the structure of experiment definitions is:

File	True Label	Train Label	Triggered?
f1	1	1	False
f2	1	2	True
...	...	...	...

Implemented `Experiment` generators are located in the `trojai.datagen.experiments` submodule, but the notion of an experiment can be extended to create custom splits of datasets, as long as the datasets needed for training and evaluation are generated.

### 3.2.1 Classic Experiment

`trojai.datagen.experiment.ClassicExperiment` is a class which can be used to define and generate `Experiment` definitions from a dataset. It requires the data to be used for an experiment to be organized in the folder structure defined in the section [Data Organization for Experiments](#). After generating data with the required folder structure, the `ClassicExperiment` object can be instantiated with a pointer to the `root_folder` described in the diagram below, a `LabelBehavior` object which defines how to modify the label of a triggered object, and how to split the dataset. Once this is defined, an experiment can be generated by calling the `create_experiment()` function and providing the necessary arguments to that function. See `trojai.datagen.experiment.ClassicExperiment` and `trojai.datagen.common_behaviors` for further details.

Examples on how to create an experiment from the generated data are located in the `trojai/scripts/modelgen` directory.

## Data Organization for Experiments

To generate experiments based on given clean data and modified data folders, the following folder structure for data is expected:

```
root_folder
|   clean_data
|       train.csv - CSV file with pointers to the training data and the associated label
+-- label
|       test.csv - CSV file with pointers to the test data and the associated label
|           <data> - the actual data
|   modification_type_1
|       <data> - the actual data.
|   modification_type_2
|       ...
|   ...
```

Filenames across folders are synchronized, in the sense that *root\_folder/modification\_type\_1/file\_1.dat* is a modified version of the file *root\_folder/clean\_data/file\_1.dat*. The same goes for *modification\_type\_2* and so on. Additionally, there are no CSV files in the modified data folders, because the required information is contained by the fact that filenames are synchronized, and the labels of those files can be referenced with the clean data CSV files.

The *train.csv* and *test.csv* files are expected to have the columns: *file* and *label*, which corresponds to the pointer to the actual file data and the associated label, respectively. Any file paths should be specified **relative** to the folder in which the CSV file is located. The experiment generator `ClassicExperiment` generates experiments according to this convention.

## 3.3 Model Generation

### 3.3.1 Overview & Concept

`trojai.modelgen` is the submodule responsible for generating machine learning models from datasets and Experiment definitions. The primary classes within `trojai.modelgen` that are of interest are:

1. `DataManager`
2. `ArchitectureFactory`
3. `OptimizerInterface`
4. `Runner`
5. `ModelGenerator`

From a top-down perspective, a `Runner` object is responsible for generating a model, trained with a given configuration specified by the `RunnerConfig`. The `RunnerConfig` consists of specifying the following parameters:

1. `ArchitectureFactory` - an object of a user-defined class which implements the interface specified by `ArchitectureFactory`. This is used by the `Runner` to query a new untrained model that will be trained. Example implementations of the `ArchitectureFactory` can be found in the scripts: `gen_and_train_mnist.py` and `gen_and_train_mnist_sequential.py`.
2. `DataManager` - an instance of the `DataManager` class, which defines the underlying datasets that will be used to train the model. Refer to the docstring for `DataManager` to understand how to instantiate this object.
3. `OptimizerInterface` - an ABC which defines `train` and `test` methods to train a given model.

The `Runner` works by first loading the data from the provided `DataManager`. Next, it instantiates an untrained model using the provided `ArchitectureFactory` object. Finally, the runner uses an optimizer specified by an

instance of an OptimizerInterface to train the model provided by the ArchitectureFactory against the data returned by the DataManager. In TrojAI nomenclature, the optimizer specifies how to train the model through the definition of the `torch.nn.module.forward()` function. Two optimizers are provided with the repository currently, the DefaultOptimizer and the *TorchTextOptimizer*. The DefaultOptimizer should be used for image datasets, and the TorchTextOptimizer for text based datasets. The RunnerConfig can accept any optimizer object that implements the OptimizerInterface, or it can accept a DefaultOptimizerConfig object and will configure the DefaultOptimizer according to the specified configuration. Thus, the Runner can be viewed a fundamental component to generate a model given a specification and corresponding configuration.

The ModelGenerator can be used to scale up model generation, by deploying the Runner in parallel on a single machine, or across a HPC cluster or cloud infrastructure. Two model generators are provided, that support single machine model generation `model_generator.py`, and HPC based model generation `uge_model_generator.py`.

### 3.3.2 Class Descriptions

#### DataManager

This object facilitates data management between the user and the module. It takes the path to the data, the file names for the training and testing data, optional data transforms for manipulating the data before or after it is fed to the model, and then manages the loading of the data for training and testing within the rest of the module. The DataManager is configured directly by the user and passed to the RunnerConfig.

#### ArchitectureFactory

This is a `factory object` which is responsible for creating new instances of trainable models. It is used by the Runner to instantiate a fresh, trainable module, to be trained by an Optimizer.

For certain model architectures or data domains, such as text, it may be the case that certain characteristics or attributes of the data are needed in order to properly setup the model that is to be trained. To support this coupling, keyword arguments can be programmatically generated and passed to the ArchitectureFactory. Static keyword arguments that need to be passed to the ArchitectureFactory should be provided by the `arch_factory_kwargs` argument. A configurable callable, which can append to the initial static arguments in `arch_factory_kwargs` can be defined via the `arch_factory_kwargs_generator` argument. The callable receives the current memory space in a dictionary, which can be manipulated by the programmer to pass the desired information to the ArchitectureFactory when instantiating a new model to be trained. Both the `arch_factory_kwargs` and `arch_factory_kwargs_generator` are optional and default to no keyword arguments being passed to the architecture factory.

Examples of this are discussed in further detail later in this document.

#### OptimizerInterface

The Runner trains a model by using a subclass of the OptimizerInterface object. The OptimizerInterface is an ABC which requires implementers to define `train` and `test` methods defining how to train and test a model. A default optimizer useful for image datasets is provided in `trojai.modelgen.default_optimizer.DefaultOptimizer`. A text optimizer is useful for text datasets and is provided in the `trojai.modelgen.torchtext_optimizer.TorchTextOptimizer`. The user is also free to specify custom training and test routines by implementing the OptimizerInterface interface.

## Adversarial Training

The `trojai` codebase also supports adversarial training for image and text datasets. Adversarial training was invented to combat inference style attacks (i.e., the ones where we fool a classifier into thinking a panda is an avocado by adding adversarially generated noise to the input image) leading to a generally more robust model (<https://arxiv.org/pdf/1412.6572.pdf>). Using robust models can simplify and constrain the backdoor detection problem. More specifically, using more robust models to train trojan detectors allows one to avoid “false positives” (ie. detecting naturally occurring triggers over the intentionally inserted triggers) and thereby study intentionally inserted triggers more effectively. Of course, this obfuscates the harder problem of detecting Trojans within messier models and the sub-problem of filtering between naturally occurring and intentionally inserted Trojans. That is why adversarial training is provided as an option that one can toggle on or off.

At a high level, adversarial training works by augmenting the training dataset with data points which are adversarially generated, while the output label is kept constant. This effectively allows one to optimize the neural network to correctly classify adversarial images, and the hope is that by training the neural network in this way, . More technically, we are seeking to minimize the empirical adversarial risk of the classifier rather than the traditional risk. For more details, visit this excellent [tutorial](#).

We have implemented two different optimizers, which implement adversarial training. The [first](#) uses the projected gradient descent (PGD) method to generate adversarial examples. Briefly, PGD generates adversarial examples by maximizing the loss of the input sample + perturbation against the output, while constraining the perturbation to be within a norm-ball. For more details, please refer to the [PGD paper](#). PGD is a general approach, but requires several iterations to find a good perturbation vector, which usually slows down training.

The second approach attempts to address the training speed by using a less computationally expensive way of generating adversarial examples, known as the fast sign gradient method (FSGM). In FSGM, adversarial examples are generated by first computing the sign the gradient of the loss with respect to the input, and then stepping the perturbation vector in the direction of the gradient by a pre-defined epsilon. This process requires only one iteration, so it is computationally much less intensive than the PGD method. However, because the iteration has no feedback, the attack is shown to be less effective, and was initially not used for adversarial training for this reason. The paper [Fast is better than free: Revisiting Adversarial Training](#) then showed that FSGM could indeed be used successfully for adversarial training if the perturbation vector is first initialized randomly and applied to the input, before computing the gradient and making a step in that direction. Here, the initialization of each element is drawn independently from the uniform distribution  $U(-\text{eps}, \text{eps})$ . This second optimizer is implemented [here](#).

## Runner

The `Runner` generates a model, given a `RunnerConfig` configuration object.

## ModelGenerator

The `ModelGenerator` is an interface for running the `Runner`, potentially parallelizing or running in parallel over a cluster or cloud interface.

For additional information about each object, see its documentation.

### 3.3.3 Scalability

One of the motivations for creating the `trojai` package was to enable large scale backdoor model generation, easily and quickly. The configuration objects and the infrastructure attempt to address the “easy” objective. To accelerate model generation, automated mixed precision (AMP) optimization is included in the `trojai` package. AMP is supported natively by PyTorch beginning with v1.7, and is effectively an engine which automatically converts some GPU operations from 32-bit floating point to 16-bit floating point (thereby increasing speed) while maintaining the same performance. It can be easily enabled when training models by simply setting the `use_amp=True` flag when configuring the `Runner` or `ModelGenerator`.

### 3.3.4 Model Generation Examples

Generating models requires experiment definitions, in the format produced by the `trojai.datagen` module. Three scripts which integrate the data generation using `trojai.datagen` submodule, and the model generation using the `trojai.modelgen` submodule are:

1. [`gen\_and\_train\_mnist.py`](#) - this script generates an MNIST dataset with an “pattern backdoor” trigger as described in the [BadNets](#) paper, and trains a model on a 20% poisoned dataset to mimic the paper’s results.
2. [`gen\_and\_train\_mnist\_sequential.py`](#) - this script generates the same MNIST dataset described above, but trains a model using an experimental feature we call “sequential” training, where the model is first trained on a clean (no-trigger) MNIST dataset and then on the poisoned dataset.
3. [`gen\_and\_train\_cifar10.py`](#) - this script generates CIFAR10 dataset with one class triggered using a Gotham Instagram filter, and trains a model on various dataset poisoning percentages.
4. [`gen\_and\_train\_imdb\_gloveilstm.py`](#) - this script generates the IMDB dataset with one class triggered with a sentence, and trains a model on various dataset poisoning percentages.



---

CHAPTER  
FOUR

---

## CONTRIBUTING

trojai welcomes your contributions! Whether it is a bug report, bug fix, new feature or documentation enhancements, please help to improve the project!

In general, please follow the [scikit-learn contribution guidelines](#) for how to contribute to an open-source project.

If you would like to open a bug report, please [open one here](#). Please try to provide a [Short, Self Contained, Example](#) so that the root cause can be pinned down and corrected more easily.

If you would like to contribute a new feature or fix an existing bug, the basic workflow to follow is:

- Open an issue with what you would like to contribute to the project and its merits. Some features may be out of scope for trojai, so be sure to get the go-ahead before working on something that is outside of the project's goals.
- Fork the trojai repository, clone it locally, and create your new feature branch.
- Make your code changes on the branch, commit them, and push to your fork.
- Open a pull request.

Please ensure that:

- Any new feature has great test coverage.
- Any new feature is well documented with [numpy-style docstrings](#) & an example, if appropriate and illustrative.
- Any bug fix has regression tests.
- Comply with [PEP8](#).



---

**CHAPTER****FIVE**

---

**ACKNOWLEDGEMENTS**

This research is based upon work supported in part by the Office of the Director of National Intelligence (ODNI), Intelligence Advanced Research Projects Activity (IARPA). The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies, either expressed or implied, of ODNI, IARPA, or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for governmental purposes notwithstanding any copyright annotation therein.



## TROJAI PACKAGE

### 6.1 Subpackages

#### 6.1.1 trojai.datagen package

##### Submodules

###### [trojai.datagen.common\\_label\\_behaviors module](#)

**class** `trojai.datagen.common_label_behaviors.StaticTarget` (`target`)

Bases: `trojai.datagen.label_behavior.LabelBehavior`

Sets label to a defined value

**do** (`y_true`)

    Performs the actual specified label modification :param `y_true`: input label to be modified :return: the modified label

**class** `trojai.datagen.common_label_behaviors.WrappedAdd` (`add_val: int, max_num_classes: int = None`)

Bases: `trojai.datagen.label_behavior.LabelBehavior`

Adds a defined amount to each input label, with an optional maximum value around which labels are wrapped

**do** (`y_true: int`) → `int`

    Performs the actual specified label modification :param `y_true`: input label to be modified :return: the modified label

`trojai.datagen.common_label_behaviors.logger = <Logger trojai.datagen.common_label_behaviors>`

Defines some common behaviors which are used to modify labels when designing an experiment with triggered and clean data

## trojai.datagen.config module

```
class trojai.datagen.config.TrojAIConfig(sign_xforms: Sequence[trojai.datagen.transform_interface.Transform] = None, bg_xforms: Sequence[trojai.datagen.transform_interface.Transform] = None, merge_obj: trojai.datagen.merge_interface.Merge = None, combined_xforms: Sequence[trojai.datagen.transform_interface.Transform] = None)
```

Bases: `object`

`validate()` → `None`

```
class trojai.datagen.config.ValidInsertLocationsConfig(algorithm: str = 'brute_force', min_val: Union[int, Sequence[int]] = 0, threshold_val: Union[float, Sequence[float]] = 5.0, num_boxes: int = 5, allow_overlap: Union[bool, Sequence[bool]] = False)
```

Bases: `object`

Specifies which algorithm to use for determining the valid spots for trigger insertion on an image and all relevant parameters

`validate()`

Assess validity of provided values :return: `None`

```
class trojai.datagen.config.XFormMergePipelineConfig(trigger_list: Sequence[trojai.datagen.entity.Entity] = None, trigger_sampling_prob: Sequence[float] = None, trigger_xforms: Sequence[trojai.datagen.transform_interface.Transform] = None, trigger_bg_xforms: Sequence[trojai.datagen.transform_interface.Transform] = None, trigger_bg_merge: trojai.datagen.merge_interface.Merge = None, trigger_bg_merge_xforms: Sequence[trojai.datagen.transform_interface.Transform] = None, over_all_bg_xforms: Sequence[trojai.datagen.transform_interface.Transform] = None, over_all_bg_triggerbg_merge: trojai.datagen.merge_interface.Merge = None, over_all_bg_triggerbg_xforms: Sequence[trojai.datagen.transform_interface.Transform] = None, merge_type: str = 'insert', per_class_trigger_frac: float = None, triggered_classes: Union[str, Sequence[Any]] = 'all')
```

Bases: `object`

Defines all configuration items necessary to run the XFormMerge Pipeline, and associated configuration validation.

NOTE: the argument list can be condensed into lists of lists, but that becomes a bit less intuitive to use. We need to think about how best we want to specify these argument lists.

#### `validate()`

Validates whether the configuration was setup properly, based on the `merge_type`. :return: `None`

#### `validate_regenerate_mode()`

Validates whether the configuration was setup properly, based on the `merge_type`. :return: `None`

`trojai.datagen.config.check_list_type(op_list, type, err_msg)`

`trojai.datagen.config.check_non_negative(val, name)`

`trojai.datagen.config.logger = <Logger trojai.datagen.config (WARNING)>`

Contains classes which define configuration used for transforming and modifying objects, as well as the associated validation routines. Ideally, a configuration class should be defined for every pipeline that is defined.

## **trojai.datagen.constants module**

`trojai.datagen.constants.RANDOM_STATE_DRAW_LIMIT = 4294967295`

In the data generation process, every new entity that is generated gets a new random seed by drawing from `np.random.RandomState.randint()`, where the `RandomState` object comes from a master `RandomState` created at the beginning of the data generation process. The constant `RANDOM_STATE_DRAW_LIMIT` defines the argument passed into the `randint(...)` call.

The reason we create a new seed for every Entity is to enable reproducibility. Each Entity that is created may go through a series of transformations that include randomness at various stages. As such, having a seed associated with each Entity will enable us to reproduce those specific random variations easily.

## **trojai.datagen.datatype\_xforms module**

`class trojai.datagen.datatype_xforms.ToTensorXForm(num_dims: int = 3)`

Bases: `trojai.datagen.transform_interface.ImageTransform`

Transformation which defines the conversion of an input array to a tensor of a specified # of dimensions

`do(input_obj: trojai.datagen.image_entity.ImageEntity, random_state_obj: numpy.random.mtrand.RandomState) → trojai.datagen.image_entity.ImageEntity`  
Perform the actual to->tensor conversion :param `input_obj`: the input Entity to be transformed :param `random_state_obj`: ignored :return: the transformed Entity

`trojai.datagen.datatype_xforms.logger = <Logger trojai.datagen.datatype_xforms (WARNING)>`

Defines data type transformations that may need to occur when processing different data sources

## **trojai.datagen.entity module**

`class trojai.datagen.entity.Entity`

Bases: `abc.ABC`

An Entity is a generalization of a synthetic object. It could stand alone, or a composition of multiple entities. An Entity is composed of some data. See the README for further details on how Entity objects are intended to be used in the TrojAI pipeline.

`abstract get_data()`

Get the data associated with the Entity :return: return the internal representation of the image

`trojai.datagen.entity.logger = <Logger trojai.datagen.entity (WARNING)>`

Defines a generic Entity object, and an Entity convenience wrapper for creating Entities from numpy arrays.

## **trojai.datagen.experiment module**

`class trojai.datagen.experiment.ClassicExperiment(data_root_dir: str, trigger_label_xform: trojai.datagen.label_behavior.LabelBehavior, stratify_split: bool = True)`

Bases: `object`

Defines a classic experiment, which consists of: 1) a specification of the clean data 2) a specification of the modified (triggered) data, and 3) a specification of the split of triggered/clean data for training/testing the model

```
create_experiment(clean_data_csv: str, experiment_data_folder: str, mod_filename_filter: str = '*', split_clean_trigger: bool = False, trigger_frac: float = 0.2, triggered_classes: Union[str, Sequence[Any]] = 'all', random_state_obj: numpy.random.mtrand.RandomState = RandomState(MT19937) at 0x7FD5BF71F5A0) → Union[Tuple, pandas.core.frame.DataFrame]
```

Creates an “experiment,” which is a dataframe defining the data that should be used, and whether that data is triggered or not, and the true & actual label associated with that data point.

**TODO:**

- [] - Have ability to accept multiple mod\_data\_folders such that we can sample from them all at a specified probability to have different triggers

### Parameters

- **clean\_data\_csv** – path to file which contains a CSV specification of the clean data. The CSV file is expected to have the following columns: [file, label]
- **experiment\_data\_folder** – the folder which contains the data to mix with for the experiment.
- **mod\_filename\_filter** – a string filter for determining which files in the folder to consider, if only a subset is to be considered for sampling
- **split\_clean\_trigger** – if True, then we return a list of DataFrames, where the triggered & non-triggered data are combined into one DataFrame, if False, we concatenate the triggered and non-triggered data into one DataFrame
- **trigger\_frac** – the fraction of data which which should be triggered
- **triggered\_classes** – either the string ‘all’, or a Sequence of labels which are to be triggered. If this parameter is ‘all’, then all classes will be triggered in the created experiment. Otherwise, only the classes in the list will be triggered at the percentage requested in the trigger\_frac argument of the create\_experiment function.
- **random\_state\_obj** – random state object

### Returns

a dataframe of the data which consists of the experiment. The DataFrame has the following columns: file, true\_label, train\_label, triggered file - the file path of the data true\_label - the actual label of the data train\_label - the label of the data the model should be trained on.

This will be equal to true\_label if triggered==False

triggered - a boolean value indicating whether this particular sample has a Trigger or not

```
trojai.datagen.experiment.logger = <Logger trojai.datagen.experiment (WARNING)>
Module which contains functionality for generating experiments
```

**trojai.datagen.image\_affine\_xforms module**

**class** trojai.datagen.image\_affine\_xforms.**PerspectiveXForm**(*xform\_matrix*)  
Bases: *trojai.datagen.transform\_interface.ImageTransform*

Shifts the perspective of an input Entity

**do** (*input\_obj*: *trojai.datagen.image\_entity.ImageEntity*, *random\_state\_obj*: *numpy.random.mtrand.RandomState*) → *trojai.datagen.image\_entity.ImageEntity*  
Performs the perspective shift on the input Entity. :param *input\_obj*: the Entity to be transformed according to the specified perspective shift in the constructor. :param *random\_state\_obj*: ignored :return: the transformed Entity

**class** trojai.datagen.image\_affine\_xforms.**RandomPerspectiveXForm**(*perspectives*: *Sequence[str]* = *None*)  
Bases: *trojai.datagen.transform\_interface.ImageTransform*

Randomly shifts perspective of input Entity in available perspectives.

**do** (*input\_obj*: *trojai.datagen.image\_entity.ImageEntity*, *random\_state\_obj*: *numpy.random.mtrand.RandomState*) → *trojai.datagen.image\_entity.ImageEntity*  
Samples from the possible perspectives according to the sampler specification and then applies that perspective to the input object :param *input\_obj*: Entity to be randomly perspective shifted :param *random\_state\_obj*: allows for reproducible sampling of random perspectives :return: the transformed Entity

**class** trojai.datagen.image\_affine\_xforms.**RandomRotateXForm**(*angle\_choices*: *Sequence[float]* = *None*, *angle\_sampler\_prob*: *Sequence[float]* = *None*, *rotator\_kwargs*: *Dict* = *None*)  
Bases: *trojai.datagen.transform\_interface.ImageTransform*

Implements a rotation of a random amount of degrees.

**do** (*input\_obj*: *trojai.datagen.image\_entity.ImageEntity*, *random\_state\_obj*: *numpy.random.mtrand.RandomState*) → *trojai.datagen.image\_entity.ImageEntity*  
Samples from the possible angles according to the sampler specification and then applies that rotation to the input object :param *input\_obj*: Entity to be randomly rotated :param *random\_state\_obj*: a random state used to maintain reproducibility through transformations :return: the transformed Entity

**class** trojai.datagen.image\_affine\_xforms.**RotateXForm**(*angle*: *int* = 90, *args*: *tuple* = (), *kwargs*: *dict* = *None*)  
Bases: *trojai.datagen.transform\_interface.ImageTransform*

Implements a rotation of an Entity by a specified angle amount.

**do** (*input\_obj*: *trojai.datagen.image\_entity.ImageEntity*, *random\_state\_obj*: *numpy.random.mtrand.RandomState*) → *trojai.datagen.image\_entity.ImageEntity*  
Performs the rotation specified by the RotateXForm object on an input :param *input\_obj*: The Entity to be rotated :param *random\_state\_obj*: ignored :return: the transformed Entity

**class** trojai.datagen.image\_affine\_xforms.**UniformScaleXForm**(*scale\_factor*: *float* = 1, *kwargs*: *dict* = *None*)  
Bases: *trojai.datagen.transform\_interface.ImageTransform*

Implements a uniform scale of a specified amount to an Entity

**do** (*input\_obj*: *trojai.datagen.image\_entity.ImageEntity*, *random\_state\_obj*: *numpy.random.mtrand.RandomState*) → *trojai.datagen.image\_entity.ImageEntity*  
Performs the scaling on an input Entity using skimage.transform.rescale :param *input\_obj*: the input object

to be scaled :param random\_state\_obj: ignored :return: the transformed Entity

```
trojai.datagen.image_affine_xforms.get_predefined_perspective_xform_matrix(xform_str:
    str,
    rows:
    int,
    cols:
    int)
    →
    numpy.ndarray
```

Returns an affine transform matrix for a string specification of a perspective transformation :param xform\_str: a string specification of the perspective to transform the object into.

### Parameters

- **rows** – the number of rows of the image to be transformed to the specified perspective
- **cols** – the number of cols of the image to be transformed to the specified perspective

**Returns** a numpy array of shape (2,3) which specifies the affine transformation.

See:[https://docs.opencv.org/2.4/modules/imgproc/doc/geometric\\_transformations.html?highlight=getaffinetransform](https://docs.opencv.org/2.4/modules/imgproc/doc/geometric_transformations.html?highlight=getaffinetransform) for more information

```
trojai.datagen.image_affine_xforms.logger = <Logger trojai.datagen.image_affine_xforms (WAI)
```

Module defines several affine transforms using various libraries to perform the actual transformation operation specified.

## trojai.datagen.image\_conversion\_utils module

```
trojai.datagen.image_conversion_utils.gray_to_rgb(img:      numpy.ndarray) →
    numpy.ndarray
Convert given grayscale image to RGB :param img: 1-channel grayscale image :return: image converted to RGB
```

```
trojai.datagen.image_conversion_utils.logger = <Logger trojai.datagen.image_conversion_utils (WAI)
Contains general utilities for dealing with channel formats
```

```
trojai.datagen.image_conversion_utils.normalization_from_rgb(rgb_img:
    numpy.ndarray,
    alpha_ch:      Optional[numpy.ndarray],
    normalize:    bool,
    original_n_chan:
    int, name: str) →
    numpy.ndarray
Guard for output from rgb-only xforms :param rgb_img: 3-channel RGB image result from calling xform :param alpha_ch: alpha channel extracted at beginning of calling xform or None :param normalize: whether to convert rgb_img back to its original channel format :param original_n_chan: number of channels in its original channel format :param name: name of calling xform :return: if normalize is True the image corresponding to rgb_img converted to its original channel format, otherwise rgb_img unmodified, additional conversions can be added below, currently only RGB to RGBA is implemented
```

```
trojai.datagen.image_conversion_utils.normalization_to_rgb(img: numpy.ndarray,  
                           normalize: bool,  
                           name: str) → Tuple[numpy.ndarray,  
                           Optional[numpy.ndarray]]
```

Guard for input to RGB only xforms :param img: input image with variable number of channels :param normalize: whether to attempt to convert img from original channel format to 3-channel RGB :param name: name of calling xform :return: a 3-channel RGB array converted from img, additional conversions can be added below, currently only RGBA to RGB is implemented

```
trojai.datagen.image_conversion_utils.rgb_to_rgba(img, alpha_ch: Optional[numpy.ndarray] = None) → numpy.ndarray
```

Converts given image to RGBA, with optionally provided alpha\_ch as its alpha channel :param img: 3-channel RGB image or 4-channel RGBA image :param alpha\_ch: 1-channel array to be used as alpha value (optional), if img is RGBA this value is ignored :return: if img is 4-channel it is returned unmodified, if img is 3-channel this will return a new RGBA image with img as its RGB channels and either alpha\_ch as its alpha channel if provided or a fully opaque alpha channel (max value for its datatype)

```
trojai.datagen.image_conversion_utils.rgba_to_rgb(img: numpy.ndarray) → Tuple[numpy.ndarray, Optional[numpy.ndarray]]
```

Split given 4-channel RGBA array into a 3-channel RGB array and a 1-channel alpha array :param img: given image to split, must be 3-channel or 4-channel :return: the first three channels of data as a 3-channel RGB image and the fourth channel of img as either a 1-channel alpha array, or None if img has only 3 channels

## trojai.datagen.image\_entity module

```
class trojai.datagen.image_entity.GenericImageEntity(data: numpy.ndarray, mask: numpy.ndarray = None)
```

Bases: *trojai.datagen.image\_entity.ImageEntity*

A class which allows one to easily instantiate an ImageEntity object with an image and associated mask

```
get_data() → numpy.ndarray
```

Get the data associated with the ImageEntity :return: return a numpy.ndarray representing the image

```
get_mask() → numpy.ndarray
```

Get the mask associated with the ImageEntity :return: return a numpy.ndarray representing the mask

```
class trojai.datagen.image_entity.ImageEntity
```

Bases: *trojai.datagen.entity.Entity*

```
abstract get_mask() → numpy.ndarray
```

```
trojai.datagen.image_entity.logger = <Logger trojai.datagen.image_entity (WARNING)>
```

Defines a generic Entity object, and an Entity convenience wrapper for creating Entities from numpy arrays.

## trojai.datagen.image\_insert\_utils module

```
trojai.datagen.image_insert_utils.pattern_fit (chan_img: numpy.ndarray, chan_pattern: numpy.ndarray, chan_location: Sequence[Any]) → bool
```

Returns True if the pattern at the desired location can fit into the image channel without wrap, and False otherwise

### Parameters

- **chan\_img** – a numpy.ndarray of shape (nrows, ncols) which represents an image channel
- **chan\_pattern** – a numpy.ndarray of shape (prows, pcols) which represents a channel of the pattern
- **chan\_location** – a Sequence of length 2, which contains the x/y coordinate of the top left corner of the pattern to be inserted for this specific channel

**Returns** True/False depending on whether the pattern will fit into the image

```
trojai.datagen.image_insert_utils.valid_locations (img: numpy.ndarray, pattern: numpy.ndarray, algo_config: trojai.datagen.config.ValidInsertLocationsConfig, protect_wrap: bool = True) → numpy.ndarray
```

Returns a list of locations per channel which the pattern can be inserted into the img\_channel with an overlap algorithm dictated by the appropriate inputs

### Parameters

- **img** – a numpy.ndarray which represents the image of shape: (nrows, ncols, nchans)
- **pattern** – the pattern to be inserted into the image of shape: (prows, pcols, nchans)
- **algo\_config** – The provided configuration object specifying the algorithm to use and necessary parameters
- **protect\_wrap** – if True, ensures that pattern to be inserted can fit without wrapping and raises an Exception otherwise

**Returns** A boolean mask of the same shape as the input image, with True indicating that that pixel is a valid location for placement of the specified pattern

## trojai.datagen.image\_size\_xforms module

```
class trojai.datagen.image_size_xforms.Pad (pad_amounts: tuple = (0, 0, 0, 0), mode: str = 'constant', pad_value: int = 0)
```

Bases: *trojai.datagen.transform\_interface.Transform*

Resizes an Entity

```
do (img_obj: trojai.datagen.image_entity.ImageEntity, random_state_obj: numpy.random.mtrand.RandomState) → trojai.datagen.image_entity.ImageEntity
```

Perform the resizing :param img\_obj: The input object to be resized according the specified configuration  
:param random\_state\_obj: ignored :return: The resized object

```
class trojai.datagen.image_size_xforms.RandomPadToSize (new_size: tuple = (200, 200), mode: str = 'constant', pad_value: int = 0)
```

Bases: *trojai.datagen.transform\_interface.Transform*

Resizes an Entity

```
do (img_obj: trojai.datagen.image_entity.ImageEntity, random_state_obj:  
    numpy.random.mtrand.RandomState) → trojai.datagen.image_entity.ImageEntity  
    Perform the resizing :param img_obj: The input object to be resized according the specified configuration  
    :param random_state_obj: ignored :return: The resized object  
  
class trojai.datagen.image_size_xforms.RandomResize (new_size_minimum: tuple =  
    (200, 200), new_size_maximum: tuple = (300, 300), interpolation:  
    int = 2)  
Bases: trojai.datagen.transform_interface.Transform  
  
Resizes an Entity  
  
do (img_obj: trojai.datagen.image_entity.ImageEntity, random_state_obj:  
    numpy.random.mtrand.RandomState) → trojai.datagen.image_entity.ImageEntity  
    Perform the resizing :param img_obj: The input object to be resized according the specified configuration  
    :param random_state_obj: ignored :return: The resized object  
  
class trojai.datagen.image_size_xforms.RandomSubCrop (new_size: tuple = (200, 200))  
Bases: trojai.datagen.transform_interface.Transform  
  
Resizes an Entity  
  
do (img_obj: trojai.datagen.image_entity.ImageEntity, random_state_obj:  
    numpy.random.mtrand.RandomState) → trojai.datagen.image_entity.ImageEntity  
    Perform the resizing :param img_obj: The input object to be cropped according the specified configuration  
    :param random_state_obj: ignored :return: The cropped object  
  
class trojai.datagen.image_size_xforms.Resize (new_size: tuple = (200, 200), interpolation:  
    int = 2)  
Bases: trojai.datagen.transform_interface.Transform  
  
Resizes an Entity  
  
do (img_obj: trojai.datagen.image_entity.ImageEntity, random_state_obj:  
    numpy.random.mtrand.RandomState) → trojai.datagen.image_entity.ImageEntity  
    Perform the resizing :param img_obj: The input object to be resized according the specified configuration  
    :param random_state_obj: ignored :return: The resized object  
  
trojai.datagen.image_size_xforms.logger = <Logger trojai.datagen.image_size_xforms> (WARNING)  
Module contains various classes that relate to size transformations of input objects
```

## `trojai.datagen.image_triggers` module

```
class trojai.datagen.image_triggers.RandomRectangularPattern (num_rows: int,  

                                                               num_cols: int,  

                                                               num_chan: int,  

                                                               color_algorithm:  

                                                               str = 'chan-  

                                                               nel_assign',  

                                                               color_options:  

                                                               dict = None, pat-  

                                                               tern_style='graffiti',  

                                                               dtype=<class  

                                                               'numpy.uint8'>,  

                                                               random_state_obj:  

                                                               numpy.random.mtrand.RandomState  

                                                               = Random-  

                                                               State(MT19937) at  

                                                               0x7FD5B287A160)
```

Bases: *trojai.datagen.image\_entity.ImageEntity*

Defines a random rectangular pattern

**create**() → None  
Create the actual pattern :return: None

**get\_data**() → numpy.ndarray  
Get the image associated with the Entity :return: return a numpy.ndarray representing the image

**get\_mask**() → numpy.ndarray  
Get the mask associated with the Entity :return: return a numpy.ndarray representing the mask

```
class trojai.datagen.image_triggers.RectangularPattern (num_rows: int, num_cols:  

                                                               int, num_chan: int,  

                                                               cval: int, dtype=<class  

                                                               'numpy.uint8'>)
```

Bases: *trojai.datagen.image\_entity.ImageEntity*

Define a rectangular pattern

**create**() → None  
Create the actual pattern :return: None

**get\_data**() → numpy.ndarray  
Get the image associated with the Entity :return: return a numpy.ndarray representing the image

**get\_mask**() → numpy.ndarray  
Get the mask associated with the Entity :return: return a numpy.ndarray representing the mask

```
class trojai.datagen.image_triggers.ReverseLambdaPattern (num_rows: int,  

                                                               num_cols: int,  

                                                               num_chan: int, trig-  

                                                               ger_cval: Union[int,  

                                                               Sequence[int]],  

                                                               bg_cval: Union[int,  

                                                               Sequence[int]] = 0,  

                                                               thickness: int = 1,  

                                                               pattern_style: str =  

                                                               'graffiti', dtype=<class  

                                                               'numpy.uint8'>)
```

Bases: *trojai.datagen.image\_entity.ImageEntity*

Defines an alpha pattern

**create()** → None

Creates the alpha pattern and associated mask :return: None

**get\_data()** → numpy.ndarray

Get the image associated with the Entity :return: return a numpy.ndarray representing the image

**get\_mask()** → numpy.ndarray

Get the mask associated with the Entity :return: return a numpy.ndarray representing the mask

```
trojai.datagen.image_triggers.logger = <Logger trojai.datagen.image_triggers (WARNING)>
Defines various Trigger Entity objects
```

## trojai.datagen.insert\_merges module

```
class trojai.datagen.insert_merges.FixedInsertTextMerge(location: int)
Bases: trojai.datagen.merge_interface.TextMerge
```

**do** (obj1: trojai.datagen.text\_entity.TextEntity, obj2: trojai.datagen.text\_entity.TextEntity, random\_state\_obj: numpy.random.mtrand.RandomState)  
Perform the actual merge operation :param obj1: the first Entity to be merged :param obj2: the second Entity to be merged :param random\_state\_obj: a numpy.random.RandomState object to ensure reproducibility :return: the merged Entity

```
class trojai.datagen.insert_merges.InsertAtLocation(location: numpy.ndarray, protect_wrap: bool = True)
Bases: trojai.datagen.merge_interface.ImageMerge
```

Inserts a provided pattern at a specified location

**do** (img\_obj: trojai.datagen.image\_entity.ImageEntity, pattern\_obj: trojai.datagen.image\_entity.ImageEntity, random\_state\_obj: numpy.random.mtrand.RandomState) → trojai.datagen.image\_entity.ImageEntity  
Inserts a pattern into an image, using the mask of the pattern to determine which specific pixels are modifiable :param img\_obj: The background image into which the pattern is inserted :param pattern\_obj: The pattern to be inserted. The mask associated with the pattern is used to determine which specific pixels of the pattern are inserted into the img\_obj

**Parameters** `random_state_obj` – ignored

**Returns** The merged object

```
class trojai.datagen.insert_merges.InsertAtRandomLocation(method: str,
algo_config: trojai.datagen.config.ValidInsertLocationsConfig,
protect_wrap: bool = True)
Bases: trojai.datagen.merge_interface.ImageMerge
```

Inserts a provided pattern at a random location, where valid locations are determined according to a provided algorithm specification

**do** (img\_obj: trojai.datagen.image\_entity.ImageEntity, pattern\_obj: trojai.datagen.image\_entity.ImageEntity, random\_state\_obj: numpy.random.mtrand.RandomState) → trojai.datagen.image\_entity.ImageEntity  
Perform the specified merge on the input Entities and return the merged Entity :param img\_obj: the image object into which the pattern is to be inserted :param pattern\_obj: the pattern object to be inserted :param random\_state\_obj: used to sample from the possible valid locations, by providing a random state,

we ensure reproducibility of the data

**Returns** the merged Entity

**class** trojai.datagen.insert\_merges.**InsertRandomLocationNonzeroAlpha**

Bases: *trojai.datagen.merge\_interface.ImageMerge*

Inserts a defined pattern into an image in a randomly selected location where the alpha channel is non-zero

**do** (*img\_obj*: trojai.datagen.image\_entity.ImageEntity, *pattern\_obj*: trojai.datagen.image\_entity.ImageEntity, *random\_state\_obj*: numpy.random.mtrand.RandomState) → trojai.datagen.image\_entity.ImageEntity

Perform the described merge operation :param *img\_obj*: The input object into which the pattern is to be inserted :param *pattern\_obj*: The pattern object which is to be inserted into the image :param *random\_state\_obj*: used to sample from the possible valid locations, by providing a random state,

we ensure reproducibility of the data

**Returns** the merged object

**class** trojai.datagen.insert\_merges.**InsertRandomWithMask**

Bases: *trojai.datagen.merge\_interface.ImageMerge*

Inserts a defined pattern into an image in a randomly selected location where the specified mask is True

**do** (*img\_obj*: trojai.datagen.image\_entity.ImageEntity, *pattern\_obj*: trojai.datagen.image\_entity.ImageEntity, *random\_state\_obj*: numpy.random.mtrand.RandomState) → trojai.datagen.image\_entity.ImageEntity

Perform the described merge operation :param *img\_obj*: The input object into which the pattern is to be inserted :param *pattern\_obj*: The pattern object which is to be inserted into the image :param *random\_state\_obj*: used to sample from the possible valid locations, by providing a random state,

we ensure reproducibility of the data

**Returns** the merged object

**class** trojai.datagen.insert\_merges.**RandomInsertTextMerge**

Bases: *trojai.datagen.merge\_interface.TextMerge*

**do** (*obj1*: trojai.datagen.text\_entity.TextEntity, *obj2*: trojai.datagen.text\_entity.TextEntity, *random\_state\_obj*: numpy.random.mtrand.RandomState)

Perform the actual merge operation :param *obj1*: the first Entity to be merged :param *obj2*: the second Entity to be merged :param *random\_state\_obj*: a numpy.random.RandomState object to ensure reproducibility :return: the merged Entity

trojai.datagen.insert\_merges.**logger = <Logger trojai.datagen.insert\_merges (WARNING)>**

Module which defines several insert style merge operations.

## trojai.datagen.instagram\_xforms module

```
class trojai.datagen.instagram_xforms.FilterXForm(channel_order: str = 'BGR',
                                                    pre_normalize: bool = True,
                                                    post_normalize: bool = True)
```

Bases: *trojai.datagen.transform\_interface.ImageTransform*

Create filter xform, if no channel order is specified it is assumed to be in BGR order (opencv default), this refers only to the first 3 channels of input data as the alpha channel is handled independently

```
do (input_obj: trojai.datagen.image_entity.ImageEntity, random_state_obj: numpy.random.mtrand.RandomState) → trojai.datagen.image_entity.ImageEntity
    Compresses 3-channel image input image as a specified filetype and stores in memory, passes to into wand and applies filter, stores filtered image as specified filetype again in memory, which is then decompressed back into 3-channel image :param input_obj: entity to be transformed :param random_state_obj: object to hold random state and enable reproducibility :return: new entity with transform applied
```

```
abstract filter(image: wand.image.Image) → wand.image.Image
    subclass defined function to be called by do :param image: wand Image to be filtered :return: filtered wand Image
```

```
class trojai.datagen.instagram_xforms.GothamFilterXForm(channel_order: str =
                                                        'BGR', pre_normalize:
                                                        bool = True,
                                                        post_normalize: bool
                                                        = True)
```

Bases: *trojai.datagen.instagram\_xforms.FilterXForm*

Class implementing Instagram's Gotham filter

```
filter(image: wand.image.Image) → wand.image.Image
    modified from https://github.com/acoomans/instagram-filters/blob/master/instagram\_filters/filters/gotham.py :param image: provided image :return: new filtered image
```

```
class trojai.datagen.instagram_xforms.KelvinFilterXForm(channel_order: str =
                                                        'BGR', pre_normalize:
                                                        bool = True,
                                                        post_normalize: bool
                                                        = True)
```

Bases: *trojai.datagen.instagram\_xforms.FilterXForm*

Class implementing Instagram's Kelvin filter

```
filter(image: wand.image.Image) → wand.image.Image
    modified from https://github.com/acoomans/instagram-filters/blob/master/instagram\_filters/filters/kelvin.py :param image: provided image :return: new filtered image
```

```
class trojai.datagen.instagram_xforms.LomoFilterXForm(channel_order: str = 'BGR',
                                                       pre_normalize: bool = True,
                                                       post_normalize: bool = True)
```

Bases: *trojai.datagen.instagram\_xforms.FilterXForm*

Class implementing Instagram's Lomo filter

```
filter(image: wand.image.Image) → wand.image.Image
    modified from https://github.com/acoomans/instagram-filters/blob/master/instagram\_filters/filters/lomo.py :param image: provided image :return: new filtered image
```

```
class trojai.datagen.instagram_xforms.NashvilleFilterXForm(channel_order: str = 'BGR', pre_normalize: bool = True, post_normalize: bool = True)
    Bases: trojai.datagen.instagram_xforms.FilterXForm
    Class implementing Instagram's Nashville filter
    filter(image: wand.image.Image) → wand.image.Image
        modified from https://github.com/acoomans/instagram-filters/blob/master/instagram\_filters/filters/nashville.py :param image: :return: new filtered image

class trojai.datagen.instagram_xforms.NoOpFilterXForm(channel_order: str = 'BGR', pre_normalize: bool = True, post_normalize: bool = True)
    Bases: trojai.datagen.instagram_xforms.FilterXForm
    No operation Transform for testing purposes
    filter(image: wand.image.Image) → wand.image.Image
        subclass defined function to be called by do :param image: wand Image to be filtered :return: filtered wand Image

class trojai.datagen.instagram_xforms.ToasterXForm(channel_order: str = 'BGR', pre_normalize: bool = True, post_normalize: bool = True)
    Bases: trojai.datagen.instagram_xforms.FilterXForm
    Class implementing Instagram's Toaster filter
    filter(image: wand.image.Image) → wand.image.Image
        modified from https://github.com/acoomans/instagram-filters/blob/master/instagram\_filters/filters/toaster.py :param image: provided image :return: new filtered image
```

## **trojai.datagen.label\_behavior module**

```
class trojai.datagen.label_behavior.LabelBehavior
    Bases: abc.ABC
    A LabelBehavior is an operation performed on the “true” label to
    abstract do(input_label: int) → int
        Perform the actual desired label manipulation :param input_label: the input label to be manipulated :return: the manipulated label
```

## **trojai.datagen.merge\_interface module**

```
class trojai.datagen.merge_interface.ImageMerge
    Bases: trojai.datagen.merge_interface.Merge
    Subclass of merges for image entities. Prevents the usage of a text merge on an image entity, which has a distinct underlying data structure.
    abstract do(obj1: trojai.datagen.image_entity.ImageEntity, obj2: trojai.datagen.image_entity.ImageEntity, random_state_obj: numpy.random.mtrand.RandomState) → trojai.datagen.image_entity.ImageEntity
        Perform the actual merge operation :param obj1: the first Entity to be merged :param obj2: the second
```

Entity to be merged :param random\_state\_obj: a numpy.random.RandomState object to ensure reproducibility :return: the merged Entity

**class** trojai.datagen.merge\_interface.**Merge**

Bases: abc.ABC

A Merge is defined as an operation on two Entities and returns a single Entity

**abstract do** (obj1: trojai.datagen.entity.Entity, obj2: trojai.datagen.entity.Entity, random\_state\_obj:

numpy.random.mtrand.RandomState) → trojai.datagen.entity.Entity

Perform the actual merge operation :param obj1: the first Entity to be merged :param obj2: the second Entity to be merged :param random\_state\_obj: a numpy.random.RandomState object to ensure reproducibility :return: the merged Entity

**class** trojai.datagen.merge\_interface.**TextMerge**

Bases: trojai.datagen.merge\_interface.Merge

Subclass of merges for text entities. Prevents the usage of an image merge on a text entity, which has a distinct underlying data structure.

**abstract do** (obj1: trojai.datagen.text\_entity.TextEntity, obj2: trojai.datagen.text\_entity.TextEntity, random\_state\_obj: numpy.random.mtrand.RandomState) → trojai.datagen.text\_entity.TextEntity

Perform the actual merge operation :param obj1: the first Entity to be merged :param obj2: the second Entity to be merged :param random\_state\_obj: a numpy.random.RandomState object to ensure reproducibility :return: the merged Entity

## trojai.datagen.pipeline module

**class** trojai.datagen.pipeline.**Pipeline**

Bases: object

A pipeline is a composition of Entities, Transforms, and Merges to produce an output Entity

**abstract process** (imglist: Iterable[trojai.datagen.entity.Entity], random\_state\_obj: numpy.random.mtrand.RandomState) → trojai.datagen.entity.Entity

The method which executes the pipeline, moving data through each of Transform & Merge objects, with data flow being defined by the implementation. :param imglist: A list of Entity objects to be processed by the Pipeline :param random\_state\_obj: a random state to pass to the transforms and merge operation to ensure

reproducibility of Entities produced by the pipeline

**Returns** The output of the pipeline

## trojai.datagen.static\_color\_xforms module

**class** trojai.datagen.static\_color\_xforms.**GrayscaleToRGBXForm**

Bases: trojai.datagen.transform\_interface.Transform

Converts an 3-channel grayscale image to RGB

**do** (input\_obj: trojai.datagen.image\_entity.ImageEntity, random\_state\_obj: numpy.random.mtrand.RandomState) → trojai.datagen.image\_entity.ImageEntity

Convert the input object from 3-channel grayscale to RGB :param input\_obj: Entity to be colorized :param random\_state\_obj: ignored :return: The colorized entity

**class** trojai.datagen.static\_color\_xforms.**RGBAtoRGB**

Bases: trojai.datagen.transform\_interface.Transform

Converts input Entity from RGBA to RGB

```
do (input_obj: trojai.datagen.image_entity.ImageEntity, random_state_obj:  

    numpy.random.mtrand.RandomState) → trojai.datagen.image_entity.ImageEntity  

    Perform the RGBA to RGB transformation :param input_obj: the Entity to be transformed :param ran-  

    dom_state_obj: ignored :return: the transformed Entity
```

```
class trojai.datagen.static_color_xforms.RGBtoRGB  
Bases: trojai.datagen.transform_interface.Transform
```

Converts input Entity from RGB to RGBA

```
do (input_obj: trojai.datagen.image_entity.ImageEntity, random_state_obj:  

    numpy.random.mtrand.RandomState) → trojai.datagen.image_entity.ImageEntity  

    Perform the RGB to RGBA transformation :param input_obj: the Entity to be transformed :param ran-  

    dom_state_obj: ignored :return: the transformed Entity
```

```
trojai.datagen.static_color_xforms.logger = <Logger trojai.datagen.static_color_xforms (WARNING)  
Defines several transformations related to static (non-random) color manipulation
```

## **trojai.datagen.transform\_interface module**

```
class trojai.datagen.transform_interface.ImageTransform  
Bases: trojai.datagen.transform_interface.Transform
```

A Transform specific to ImageEntity objects

```
abstract do (input_obj: trojai.datagen.image_entity.ImageEntity, random_state_obj:  

    numpy.random.mtrand.RandomState) → trojai.datagen.image_entity.ImageEntity  

    Perform the specified transformation :param input_obj: the input ImageEntity to be transformed :param ran-  

    dom_state_obj: a random state used to maintain reproducibility through transformations :return: the  

    transformed ImageEntity
```

```
class trojai.datagen.transform_interface.TextTransform  
Bases: trojai.datagen.transform_interface.Transform
```

A Transform specific to TextEntity objects

```
abstract do (input_obj: trojai.datagen.text_entity.TextEntity, random_state_obj:  

    numpy.random.mtrand.RandomState) → trojai.datagen.text_entity.TextEntity  

    Perform the specified transformation :param input_obj: the input TextEntity to be transformed :param ran-  

    dom_state_obj: a random state used to maintain reproducibility through transformations :return: the  

    transformed TextEntity
```

```
class trojai.datagen.transform_interface.Transform
```

Bases: *abc.ABC*

A Transform is defined as an operation on an Entity.

```
abstract do (input_obj: trojai.datagen.entity.Entity, random_state_obj:  

    numpy.random.mtrand.RandomState) → trojai.datagen.entity.Entity  

    Perform the specified transformation :param input_obj: the input Entity to be transformed :param ran-  

    dom_state_obj: a random state used to maintain reproducibility through transformations :return: the trans-  

    formed Entity
```

## trojai.datagen.utils module

```
trojai.datagen.utils.logger = <Logger trojai.datagen.utils (WARNING)>
```

Contains general utilities helpful for data generation

```
trojai.datagen.utils.process_xform_list(input_obj: trojai.datagen.entity.Entity, xforms: Iterable[trojai.datagen.transform_interface.Transform], random_state_obj: numpy.random.mtrand.RandomState) → trojai.datagen.entity.Entity
```

Processes a list of transformations in a serial fashion on a copy of the input X :param input\_obj: input object which should be transformed by the list of

transformations

### Parameters

- **xforms** – a list of Transform objects
- **random\_state\_obj** –

**Returns** The transformed object

## trojai.datagen.xform\_merge\_pipeline module

```
class trojai.datagen.xform_merge_pipeline.XFormMerge(xform_list: Sequence[Sequence[Sequence[trojai.datagen.transform_interface.Transform]]], merge_list: Sequence[trojai.datagen.merge_interface.Merge], final_xforms: Sequence[trojai.datagen.transform_interface.Transform] = None)
```

Bases: *trojai.datagen.pipeline.Pipeline*

Implements a pipeline which is a series of cascading transform and merge operations. The following diagram shows 4 objects as a series of serial transforms + merges. Each pair of transformations is considered a “stage”, and stages are processed in serial fashion. In the diagram below, the data that each stage processes is:

Stage1: obj1, obj2 Stage2: Stage1\_output, obj3 Stage3: Stage2\_output, obj4

This extends in the obvious way to more objects, depending on how deep the pipeline is.

**obj1 → xform obj3 → xform obj4 → xform + → xform → + → xform → + → xform output /**

obj2 → xform

```
process(imglist: Sequence[trojai.datagen.entity.Entity], random_state_obj: numpy.random.mtrand.RandomState) → trojai.datagen.entity.Entity
```

Processes the provided objects according to the Xform->Merge->Xform paradigm. :param imglist: a sequence of Entity objects to be processed according to the pipeline :param random\_state\_obj: a random state to pass to the transforms and merge operation to ensure

reproducibility of Entities produced by the pipeline

**Returns** the modified & combined Entity object

```
trojai.datagen.xform_merge_pipeline.logger = <Logger trojai.datagen.xform_merge_pipeline (WARNING)>
```

Defines all functions and classes related to the transform+merge pipeline & data movement paradigm.

```

trojai.datagen.xform_merge_pipeline.modify_clean_image_dataset(clean_dataset_rootdir:
    str,
    clean_csv_file:
    str,          out-
    put_rootdir:
    str,          out-
    put_subdir: str,
    mod_cfg: tro-
    jai.datagen.config.XFormMergePipelineC
method: str =
'insert', ran-
dom_state_obj:
numpy.random.mtrand.RandomState
= Random-
State(MT19937)
at
0x7FD5B25C4AF0)
→ None

```

Modifies a clean dataset given a configuration

#### Parameters

- **clean\_dataset\_rootdir** – root directory where the clean data lives
- **clean\_csv\_file** – filename of the CSV file which contains information about the clean data. The modification method determines which columns and information are expected in the CSV file.
- **output\_rootdir** – root directory where the modified data will be stored
- **output\_subdir** – subdirectory where the modified data will be stored. This is expected to be one level below the root-directory, and can prove useful if different types of modifications are stored in different subdirectories under the main root directory. An example tree structure might be: root\_data
  - **modification\_1** ... data ...
  - **modification\_2** ... data ...
- **mod\_cfg** – A configuration object for creating a modified dataset
- **method** – Can be “insert” only! In the insert method, the function takes the clean image, and inserts a specified Entity (likely, a pattern) into the clean image. Additional modes to be added!
- **random\_state\_obj** – RandomState object to ensure reproducibility of dataset

**Returns** None

```
trojai.datagen.xform_merge_pipeline.modify_clean_text_dataset(clean_dataset_rootdir:  
str,  
clean_csv_file:  
str,          out-  
put_rootdir: str,  
output_subdir:  
str, mod_cfg: tro-  
jai.datagen.config.XFormMergePipelineCo  
method='insert',  
ran-  
dom_state_obj:  
numpy.random.mtrand.RandomState  
=      Random-  
State(MT19937)  
at  
0x7FD5B25C4C00)  
→ None
```

Modifies a clean image dataset given a configuration

#### Parameters

- **clean\_dataset\_rootdir** – root directory where the clean data lives
- **clean\_csv\_file** – filename of the CSV file which contains information about the clean data. The modification method determines which columns and information are expected in the CSV file.
- **output\_rootdir** – root directory where the modified data will be stored
- **output\_subdir** – subdirectory where the modified data will be stored. This is expected to be one level below the root-directory, and can prove useful if different types of modifications are stored in different subdirectories under the main root directory. An example tree structure might be: root\_data
  - **modification\_1** ... data ...
  - **modification\_2** ... data ...
- **mod\_cfg** – A configuration object for creating a modified dataset
- **method** – Can only be “insert” In the insert method, the function takes the clean text blurb, and inserts a specified TextEntity (likely, a pattern) into the first text input object.
- **random\_state\_obj** – RandomState object to ensure reproduciblity of dataset

#### Returns

```
trojai.datagen.xform_merge_pipeline.subset_clean_df_by_labels(df,           la-  
bel_to_include)
```

Subsets a dataframe with an expected column ‘label’, to only keep rows which are in that list of labels to include :param df: the dataframe to subset :param labels\_to\_include: a list of labels to include, or a string ‘all’ indicating that everything should be kept :return: the subsetted data frame

## Module contents

### 6.1.2 trojai.modelgen package

#### Subpackages

##### trojai.modelgen.architectures package

#### Submodules

##### trojai.modelgen.architectures.cifar10\_architectures module

```
class trojai.modelgen.architectures.cifar10_architectures.AlexNet (num_classes=10)
Bases: torch.nn.Module

Modified AlexNet for CIFAR From: https://github.com/icpm/pytorch-cifar10/blob/master/models/AlexNet.py

forward(x)

class trojai.modelgen.architectures.cifar10_architectures.Bottleneck (in_planes,
                                                                     growth_rate)
Bases: torch.nn.Module

Bottleneck module in DenseNet Arch. See: https://arxiv.org/abs/1608.06993

forward(x)

class trojai.modelgen.architectures.cifar10_architectures.DenseNet (block,
                                                                    num_block,
                                                                    growth_rate=12,
                                                                    reduc-
                                                                    tion=0.5,
                                                                    num_classes=10)
Bases: torch.nn.Module

From: https://github.com/icpm/pytorch-cifar10/blob/master/models/DenseNet.py

forward(x)

trojai.modelgen.architectures.cifar10_architectures.DenseNet121()
trojai.modelgen.architectures.cifar10_architectures.DenseNet161()
trojai.modelgen.architectures.cifar10_architectures.DenseNet169()
trojai.modelgen.architectures.cifar10_architectures.DenseNet201()

class trojai.modelgen.architectures.cifar10_architectures.Transition (in_planes,
                                                                     out_planes)
Bases: torch.nn.Module

Transition module in DenseNet Arch. See: https://arxiv.org/abs/1608.06993

forward(x)

trojai.modelgen.architectures.cifar10_architectures.densenet_cifar()
```

## trojai.modelgen.architectures.mnist\_architectures module

```
class trojai.modelgen.architectures.mnist_architectures.BadNetExample
Bases: torch.nn.Module

Mnist network from BadNets paper Input - 1x28x28 C1 - 1x28x28 (5x5 kernel) -> 16x24x24 ReLU S2 - 16x24x24 (2x2 kernel, stride 2) Subsampling -> 16x12x12 C3 - 16x12x12 (5x5 kernel) -> 32x8x8 ReLU S4 - 32x8x8 (2x2 kernel, stride 2) Subsampling -> 32x4x4 F6 - 512 -> 512 tanh F7 - 512 -> 10 Softmax (Output)

forward(img)

class trojai.modelgen.architectures.mnist_architectures.ModdedLeNet5Net (channels=1)
Bases: torch.nn.Module

A modified LeNet architecture that seems to be easier to embed backdoors in than the network from the original badnets paper Input - (1 or 3)x28x28 C1 - 6@28x28 (5x5 kernel) ReLU S2 - 6@14x14 (2x2 kernel, stride 2) Subsampling C3 - 16@10x10 (5x5 kernel) ReLU S4 - 16@5x5 (2x2 kernel, stride 2) Subsampling C5 - 120@1x1 (5x5 kernel) F6 - 84 ReLU F7 - 10 (Output)

forward(img)
```

### Module contents

#### Submodules

##### trojai.modelgen.architecture\_factory module

```
class trojai.modelgen.architecture_factory.ArchitectureFactory
Bases: abc.ABC

Factory object that returns architectures (untrained models) for training.

abstract new_architecture(**kwargs) → torch.nn.Module
>Returns a new architecture (untrained model) :return: an untrained torch.nn.Module
```

##### trojai.modelgen.config module

```
class trojai.modelgen.config.ConfigInterface
Bases: abc.ABC

Defines the interface for all configuration objects

class trojai.modelgen.config.DefaultOptimizerConfig (training_cfg: tro-
                                              jai.modelgen.config.TrainingConfig
                                              = None, reporting_cfg: tro-
                                              jai.modelgen.config.ReportingConfig
                                              = None)
Bases: trojai.modelgen.config.OptimizerConfigInterface

Defines the configuration needed to setup the DefaultOptimizer

get_device_type()
>Returns the device associated w/ this optimizer configuration. Needed to save/load for UGE. :return (str):
    the device type represented as a string

static load(fname)
>Loads a configuration from disk :param fname: the filename where the config is stored :return: the loaded
    configuration
```

```

save (fname)
    Saves the optimizer configuration to a file :param fname: the filename to save the config to :return: None

class trojai.modelgen.config.DefaultSoftToHardFn
    Bases: object

    The default conversion from soft-decision outputs to hard-decision

class trojai.modelgen.config.EarlyStoppingConfig (num_epochs: int = 5, val_loss_eps:
                                                    float = 0.001)
    Bases: trojai.modelgen.config.ConfigInterface

    Defines configuration related to early stopping.

validate()

class trojai.modelgen.config.ModelGeneratorConfig (arch_factory:
                                                       tro-
                                                       jai.modelgen.architecture_factory.ArchitectureFactory,
                                                       data:
                                                       tro-
                                                       jai.modelgen.data_manager.DataManager,
                                                       model_save_dir:
                                                       str,
                                                       stats_save_dir:
                                                       str,
                                                       num_models:
                                                       int,
                                                       arch_factory_kwargs: dict = None,
                                                       arch_factory_kwargs_generator:
                                                       Callable = None, optimizer:
                                                       Union[trojai.modelgen.optimizer_interface.OptimizerInterface,
                                                       tro-
                                                       jai.modelgen.config.DefaultOptimizerConfig,
                                                       Sequence[Union[trojai.modelgen.optimizer_interface.Optim-
                                                       tro-
                                                       jai.modelgen.config.DefaultOptimizerConfig]]]
                                                       = None, parallel=False,
                                                       amp=False, experiment_cfg:
                                                       dict = None, run_ids: Union[Any,
                                                       Sequence[Any]] = None, filenames:
                                                       Union[str, Sequence[str]] = None,
                                                       save_with_hash: bool = False)
    Bases: trojai.modelgen.config.ConfigInterface

    Object used to configure the model generator

static load (fname: str)
    Loads a saved modelgen_cfg object from data that was saved using the .save() function. :param fname:
        the filename where the modelgen_cfg object is saved :return: a ModelGeneratorConfig object

save (fname: str)
    Saves the ModelGeneratorConfig object in two different parts. Every object within the config, except for
    the optimizer is saved in the .klass.save file, and the optimizer is saved separately. :param fname - the
    filename to save the configuration to :return: None

validate() → None
    Validate the input arguments to construct the object :return: None

class trojai.modelgen.config.OptimizerConfigInterface
    Bases: trojai.modelgen.config.ConfigInterface

abstract get_device_type()

abstract static load (fname)

```

**save** (*fname*)

```
class trojai.modelgen.config.ReportingConfig(num_batches_per_logmsg: int = 100,
                                              disable_progress_bar: bool = False,
                                              num_epochs_per_metric: int = 1,
                                              num_batches_per_metrics: int = 50,
                                              tensorboard_output_dir: str = None,
                                              experiment_name: str = 'experiment')
```

Bases: *trojai.modelgen.config.ConfigInterface*

Defines all options to setup how data is reported back to the user while models are being trained

**validate**()

```
class trojai.modelgen.config.RunnerConfig(arch_factory: tro-
                                         jai.modelgen.architecture_factory.ArchitectureFactory,
                                         data: trojai.modelgen.data_manager.DataManager,
                                         arch_factory_kwargs: dict = None,
                                         arch_factory_kwargs_generator:
                                         Callable = None, optimizer:
                                         Union[trojai.modelgen.optimizer_interface.OptimizerInterface,
                                         trojai.modelgen.config.DefaultOptimizerConfig,
                                         Sequence[Union[trojai.modelgen.optimizer_interface.OptimizerInterface,
                                         trojai.modelgen.config.DefaultOptimizerConfig]]]
                                         = None, parallel: bool = False, amp: bool =
                                         False, model_save_dir: str = '/tmp/models',
                                         stats_save_dir: str = '/tmp/model_stats',
                                         model_save_format: str = 'pt', run_id: Any =
                                         None, filename: str = None, save_with_hash:
                                         bool = False)
```

Bases: *trojai.modelgen.config.ConfigInterface*

Container for all parameters needed to use the Runner to train a model.

**static setup\_optimizer\_generator** (*optimizer*, *data*)

Converts an optimizer specification to a generator, to be compatible with sequential training. :param *optimizer*: the optimizer to configure into a generator :param *num\_datasets*: the number of datasets for which optimizers need to be created :return: A generator that returns optimizers for every dataset to be trained

**validate**() → None

Validate the RunnerConfig object :return: None

**static validate\_optimizer** (*optimizer*, *data*)

Validates an optimizer configuration :param *optimizer*: the optimizer/optimizer configuration to be validated :param *data*: the data to be optimized :return:

```
class trojai.modelgen.config.TorchTextOptimizerConfig(training_cfg: tro-
                                                       jai.modelgen.config.TrainingConfig
                                                       = None, reporting_cfg: tro-
                                                       jai.modelgen.config.ReportingConfig
                                                       = None,
                                                       copy_pretrained_embeddings:
                                                       bool = False)
```

Bases: *trojai.modelgen.config.OptimizerConfigInterface*

Defines the configuration needed to setup the TorchTextOptimizer

**get\_device\_type**()

Returns the device associated w/ this optimizer configuration. Needed to save/load for UGE. :return (str):

the device type represented as a string

**static load**(*fname*)  
 Loads a configuration from disk :param *fname*: the filename where the config is stored :return: the loaded configuration

**save**(*fname*)  
 Saves the optimizer configuration to a file :param *fname*: the filename to save the config to :return: None

**validate**()

```
class trojai.modelgen.config.TrainingConfig(device: Union[str, torch.device] = 'cpu',
                                             epochs: int = 10, batch_size: int = 32,
                                             lr: float = 0.0001, optim: Union[str, trojai.modelgen.optimizer_interface.OptimizerInterface] = 'adam', optim_kwarg: dict = None,
                                             objective: Union[str, Callable] = 'cross_entropy_loss', objective_kwarg: dict = None, save_best_model: bool = False, train_val_split: float = 0.05,
                                             val_data_transform: Callable[[Any], Any] = None, val_label_transform: Callable[[int], int] = None, val_dataloader_kwarg: dict = None, early_stopping: trojai.modelgen.config.EarlyStoppingConfig = None, soft_to_hard_fn: Callable = None, soft_to_hard_fn_kwarg: dict = None, lr_scheduler: Any = None, lr_scheduler_init_kwarg: dict = None, lr_scheduler_call_arg: Any = None, clip_grad: bool = False, clip_type: str = 'norm', clip_val: float = 1.0, clip_kwarg: dict = None, adv_training_eps: float = None, adv_training_iterations: int = None, adv_training_ratio: float = None)
```

Bases: *trojai.modelgen.config.ConfigInterface*  
 Defines all required items to setup training with an optimizer

**get\_cfg\_as\_dict**()  
 Returns a dictionary representation of the configuration :return: (dict) a dictionary

**validate**() → None  
 Validate the object configuration :return: None

```
class trojai.modelgen.config.UGEConfig(queues: Union[trojai.modelgen.config.UGEQueueConfig, Sequence[trojai.modelgen.config.UGEQueueConfig]], queue_distribution: Sequence[float] = None, multi_model_same_gpu: bool = False)
```

Bases: *object*  
 Defines a configuration for the UGE

**validate**()  
 Validate the UGEConfig object

```
class trojai.modelgen.config.UGEQueueConfig(queue_name: str, gpu_enabled: bool, sync_mode: bool = False)
```

Bases: *object*

Defines the configuration for a Queue w.r.t. UGE in TrojAI

**validate()** → None

Validate the UGEQueueConfig object

`trojai.modelgen.config.identity_function(x)`

`trojai.modelgen.config.logger = <Logger trojai.modelgen.config (WARNING)>`

Defines all configurations pertinent to model generation.

`trojai.modelgen.config.modelgen_cfg_to_runner_cfg(modelgen_cfg: trojai.modelgen.config.ModelGeneratorConfig, run_id=None, file_name=None) → trojai.modelgen.config.RunnerConfig`

Convenience function which creates a RunnerConfig object, from a ModelGeneratorConfig object. :param modelgen\_cfg: the ModelGeneratorConfig to convert :param run\_id: run\_id to be associated with the RunnerConfig :param filename: filename to be associated with the RunnerConfig :return: the created RunnerConfig object

## **trojai.modelgen.constants module**

Defines valid devices on which models can be trained

`trojai.modelgen.constants.VALID_DEVICES = ['cpu', 'cuda']`

Defines valid loss functions which can be specified when configuring an optimizer implementing the OptimizerInterface

`trojai.modelgen.constants.VALID_LOSS_FUNCTIONS = ['cross_entropy_loss', 'BCEWithLogitsLoss']`

Defines valid optimization algorithms which can be specified when configuring an optimizer implementing the OptimizerInterface

`trojai.modelgen.constants.VALID_OPTIMIZERS = ['adam', 'sgd', 'adamm']`

Defines the valid types of data that the modelgen pipeline can handle

## **trojai.modelgen.data\_configuration module**

**class** `trojai.modelgen.data_configuration.DataConfiguration`

Bases: `object`

**class** `trojai.modelgen.data_configuration.ImageDataConfiguration`

Bases: `trojai.modelgen.data_configuration.DataConfiguration`

**class** `trojai.modelgen.data_configuration.TextDataConfiguration(max_vocab_size:`

`int = 25000, embedding_dim: int = 100, embedding_type: str = 'glove', num_tokens_embedding_train: str = '6B', text_field_kwargs: dict = None, label_field_kwargs: dict = None)`

Bases: `trojai.modelgen.data_configuration.DataConfiguration`

`set_embedding_vectors_cfg()`

```
validate()
trojai.modelgen.data_configuration.logger = <Logger trojai.modelgen.data_configuration (WARNING)
Configurations for various types of data
```

## trojai.modelgen.data\_descriptions module

File describes data description classes, which contain specific information that may be used in order to instantiate an architecture

```
class trojai.modelgen.data_descriptions.CSVImageDatasetDesc (num_samples, shuffled, num_classes)
```

Bases: *trojai.modelgen.data\_descriptions.DataDescription*

Information potentially relevant to instantiating models to process image data

```
class trojai.modelgen.data_descriptions.CSVTextDatasetDesc (vocab_size, unk_idx, pad_idx)
```

Bases: *trojai.modelgen.data\_descriptions.DataDescription*

Information potentially relevant to instantiating models to process text data

```
class trojai.modelgen.data_descriptions.DataDescription
```

Bases: *object*

Generic Data Description class from which all specific data type data descriptors

## trojai.modelgen.data\_manager module

```
class trojai.modelgen.data_managerDataManager (experiment_path: str, train_file: Union[str, Sequence[str]], clean_test_file: str, triggered_test_file: str = None, data_type: str = 'image', train_data_transform: Callable[[Any], Any] = <function DataManager.<lambda>>, train_label_transform: Callable[[int], int] = <function DataManager.<lambda>>, test_data_transform: Callable[[Any], Any] = <function DataManager.<lambda>>, test_label_transform: Callable[[int], int] = <function DataManager.<lambda>>, file_loader: Union[Callable[[str], Any], str] = 'default_image_loader', shuffle_train=True, shuffle_clean_test=False, shuffle_triggered_test=False, data_configuration: trojai.modelgen.data_configuration.DataConfiguration = None, custom_datasets: dict = None, train_dataloader_kwarg: dict = None, test_dataloader_kwarg: dict = None)
```

Bases: *object*

Manages data from an experiment from trojai.datagen.

**load\_data()**

Load experiment data as given from initialization. :return: Objects containing training and test, and triggered data if it was provided.

**TODO:** [ ] - extend the text data-type to have more input arguments, for example the tokenizer and FIELD options [ ] - need to support sequential training for text datasets

**validate() → None**

Validate the construction of the TrojaiDataManager object :return: None

**TODO:**

[ ] - think about whether the contents of the files passed into the DataManager should be validated, in addition to simply checking for existence, which is what is done now

## trojai.modelgen.datasets module

```
class trojai.modelgen.datasets.CSVDataset(path_to_data: str, csv_filename: str,
                                           true_label=False, path_to_csv=None,
                                           shuffle=False, random_state: Union[int,
                                           numpy.random.mtrand.RandomState]
                                           = None, data_loader: Union[str,
                                           Callable] = 'default_image_loader',
                                           data_transform=<function identity_transform>, label_transform=<function identity_transform>)
```

Bases: *trojai.modelgen.datasets.DatasetInterface*

Defines a dataset that is represented by a CSV file with columns “file”, “train\_label”, and optionally “true\_label”. The file column should contain the path to the file that contains the actual data, and “train\_label” refers to the label with which the data should be trained. “true\_label” refers to the actual label of the data point, and can differ from train\_label if the dataset is poisoned. A CSVDataset can support any underlying data that can be loaded on the fly and fed into the model (for example: image data)

**get\_data\_description()****set\_data\_description()**

```
class trojai.modelgen.datasets.CSVTextDataset(path_to_data: str, csv_filename: str,
                                              true_label: bool = False, text_field: torchtext.data.Field = None, text_field_kwargs: dict = None, label_field: torchtext.data.LabelField = None, label_field_kwargs: dict = None, shuffle: bool = False, random_state=None,
                                              **kwargs)
```

Bases: *torchtext.data.Dataset*, *trojai.modelgen.datasets.DatasetInterface*

Defines a text dataset that is represented by a CSV file with columns “file”, “train\_label”, and optionally “true\_label”. The file column should contain the path to the file that contains the actual data, and “train\_label” refers to the label with which the data should be trained. “true\_label” refers to the actual label of the data point, and can differ from train\_label if the dataset is poisoned. A CSVTextDataset can support text data, and differs from the CSVDataset because it loads all the text data into memory and builds a vocabulary from it.

**build\_vocab(embedding\_vectors\_cfg, max\_vocab\_size, use\_vocab=True)****get\_data\_description()****set\_data\_description()**

```

static sort_key(ex)

class trojai.modelgen.datasets.DatasetInterface(path_to_data: str, *args, **kwargs)
    Bases: torch.utils.data.Dataset

        abstract get_data_description()

        abstract set_data_description()

trojai.modelgen.datasets.csv_dataset_from_df(path_to_data, data_df, true_label=False,
                                             shuffle=False, random_state: Union[int,
                                             numpy.random.mtrand.RandomState]
                                             = None, data_loader: Union[str,
                                             Callable] = 'default_image_loader',
                                             data_transform=<function
                                             identity_transform>, label_transform=<function
                                             identity_transform>)

```

Initializes a CSVDataset object from a DataFrame rather than a filepath. :param path\_to\_data: root folder where all the data is located :param data\_df: the dataframe in which the data lives :param true\_label: (bool) if True, then use the column “true\_label” as the label associated with each datapoint. If False (default), use the column “train\_label” as the label associated with each datapoint :param shuffle: if True, the dataset is shuffled before loading into the model :param random\_state: if specified, seeds the random sampler when shuffling the data :param data\_loader: either a string value (currently only supports *default\_image\_loader*), or a callable

function which takes a string input of the file path and returns the data

#### Parameters

- **data\_transform** – a callable function which is applied to every data point before it is fed into the model. By default, this is an identity operation
- **label\_transform** – a callable function which is applied to every label before it is fed into the model. By default, this is an identity operation.

```

trojai.modelgen.datasets.csv_textdataset_from_df(data_df, true_label: bool = False,
                                              text_field: torchtext.data.Field
                                              = None, label_field: torchtext.data.LabelField = None, shuffle:
                                              bool = False, random_state=None,
                                              **kwargs)

```

Initializes a CSVDataset object from a DataFrame rather than a filepath. :param data\_df: the dataframe in which the data lives :param true\_label: if True, then use the column “true\_label” as the label associated with each :param text\_field: defines how the text data will be converted to

a Tensor. If none, a default will be provided and tokenized with spacy

#### Parameters

- **label\_field** – defines how to process the label associated with the text
- **max\_vocab\_size** – the maximum vocabulary size that will be built
- **shuffle** – if True, the dataset is shuffled before loading into the model
- **random\_state** – if specified, seeds the random sampler when shuffling the data
- **kwargs** – any additional keyword arguments, currently unused

```
trojai.modelgen.datasets.default_image_file_loader(img_loc)
```

```
trojai.modelgen.datasets.identity_transform(x)
trojai.modelgen.datasets.logger = <Logger trojai.modelgen.datasets (WARNING)>
    Define some basic default functions for dataset defaults. These allow Dataset objects to be pickled; vs lambda
    functions.
```

## trojai.modelgen.default\_optimizer module

```
class trojai.modelgen.default_optimizer.DefaultOptimizer(optimizer_cfg: tro-
    jai.modelgen.config.DefaultOptimizerConfig
    = None)
Bases: trojai.modelgen.optimizer_interface.OptimizerInterface
```

Defines the default optimizer which trains the models

**get\_cfg\_as\_dict()** → dict

Return a dictionary with key/value pairs that describe the parameters used to train the model.

**get\_device\_type()** → str

**Returns** a string representing the device used to train the model

**static load(fname: str)** → trojai.modelgen.optimizer\_interface.OptimizerInterface

Reconstructs a DefaultOptimizer, by loading the configuration used to construct the original DefaultOptimizer, and then creating a new DefaultOptimizer object from the saved configuration :param fname: The filename of the saved optimzier :return: a DefaultOptimizer object

**save(fname: str)** → None

Saves the configuration object used to construct the DefaultOptimizer. NOTE: because the DefaultOptimizer object itself is not persisted, but rather the

DefaultOptimizerConfig object, the state of the object is not persisted!

**Parameters** **fname** – the filename to save the DefaultOptimizer's configuration.

**Returns** None

```
test(net: torch.nn.Module, clean_data: trojai.modelgen.datasets.CSVDataset, trig-
    gered_data: trojai.modelgen.datasets.CSVDataset, clean_test_triggered_labels_data: tro-
    jai.modelgen.datasets.CSVDataset, torch_dataloader_kwargs: dict = None) → dict
```

Test the trained network :param net: the trained module to run the test data through :param clean\_data: the clean Dataset :param triggered\_data: the triggered Dataset, if None, not computed :param clean\_test\_triggered\_labels\_data: triggered part of the training dataset but with correct labels; see

DataManger.load\_data for more information.

**Parameters** **torch\_dataloader\_kwargs** – any keyword arguments to pass directly to PyTorch's DataLoader

**Returns** a dictionary of the statistics on the clean and triggered data (if applicable)

```
train(net: torch.nn.Module, dataset: trojai.modelgen.datasets.CSVDataset,
    torch_dataloader_kwargs: dict = None, use_amp: bool = False) -> (torch.nn.Module,
    typing.Sequence[trojai.modelgen.training_statistics.EpochStatistics], <class 'int'>)
```

Train the network. :param net: the network to train :param dataset: the dataset to train the network on :param torch\_dataloader\_kwargs: any additional kwargs to pass to PyTorch's native DataLoader :param use\_amp: if True, uses automated mixed precision for FP16 training. :return: the trained network, and a list of EpochStatistics objects which contain the statistics for training,

and the # of epochs on which the net was trained

---

**train\_epoch** (*model: torch.nn.Module*, *train\_loader: torch.utils.data.DataLoader*, *val\_clean\_loader: torch.utils.data.DataLoader*, *val\_triggered\_loader: torch.utils.data.DataLoader*, *epoch\_num: int*, *use\_amp: bool = False*)  
Runs one epoch of training on the specified model

**Parameters**

- **model** – the model to train for one epoch
- **train\_loader** – a DataLoader object pointing to the training dataset
- **val\_clean\_loader** – a DataLoader object pointing to the validation dataset that is clean
- **val\_triggered\_loader** – a DataLoader object pointing to the validation dataset that is triggered
- **epoch\_num** – the epoch number that is being trained
- **use\_amp** – if True use automated mixed precision for FP16 training.

**Returns** a list of statistics for batches where statistics were computedtrojai.modelgen.default\_optimizer.**split\_val\_clean\_trig** (*val\_dataset*)

Splits the validation dataset into clean and triggered.

**Parameters** **val\_dataset** – the validation dataset to split**Returns** A tuple of the clean & triggered validation dataset

trojai.modelgen.default\_optimizer.**train\_val\_dataset\_split** (*dataset: torch.utils.data.Dataset*, *split\_amt: float*, *val\_data\_transform: Callable*, *val\_label\_transform: Callable*) -> (*torch.utils.data.Dataset*, *torch.utils.data.Dataset*)

Splits a PyTorch dataset (of type: *torch.utils.data.Dataset*) into train/test TODO:

[ ] - specify random seed to torch splitter

**Parameters**

- **dataset** – the dataset to be split
- **split\_amt** – fraction specifying the validation dataset size relative to the whole. 1-split\_amt will be the size of the training dataset
- **val\_data\_transform** – (function: any -> any) how to transform the validation data to fit into the desired model and objective function
- **val\_label\_transform** – (function: any -> any) how to transform the validation labels

**Returns** a tuple of the train and validation datasets

## trojai.modelgen.model\_generator module

```
class trojai.modelgen.model_generator.ModelGenerator(configs:  
                                                    Union[trojai.modelgen.config.ModelGeneratorConfig,  
                                                    Sequence[trojai.modelgen.config.ModelGeneratorConfig]]  
                                                    *args, **kwargs)  
Bases: trojai.modelgen.model_generator_interface.ModelGeneratorInterface  
  
Generates models based on requested data and saves each to a file.  
  
run (*args, **kwargs) → None  
    Train and save models as specified. :return: None  
  
validate () → None  
    Validate the provided input when constructing the ModelGenerator interface
```

## trojai.modelgen.model\_generator\_interface module

```
class trojai.modelgen.model_generator_interface.ModelGeneratorInterface(configs:  
                                                    Union[trojai.modelgen.config.ModelGeneratorConfig,  
                                                    Sequence[trojai.modelgen.config.ModelGeneratorConfig]]  
                                                    *args, **kwargs)  
Bases: abc.ABC  
  
Generates models based on requested data and saves each to a file.  
  
abstract run () → None  
    Train and save models as specified. :return: None  
  
trojai.modelgen.model_generator_interface.validate_model_generator_interface_input(configs:  
                                                    Union[trojai.modelgen.config.ModelGeneratorConfig,  
                                                    Sequence[trojai.modelgen.config.ModelGeneratorConfig]]  
                                                    *args, **kwargs) → None  
  
    Validates a ModelGeneratorConfig :param configs: (ModelGeneratorConfig or sequence) configurations to be  
    used for model generation :return None
```

## trojai.modelgen.optimizer\_interface module

```
class trojai.modelgen.optimizer_interface.OptimizerInterface  
Bases: abc.ABC  
  
Object that performs training and testing of TrojAI models.  
  
abstract get_cfg_as_dict () → dict  
    Return a dictionary with key/value pairs that describe the parameters used to train the model.  
  
abstract get_device_type () → str  
    Return a string representation of the type of device used by the optimizer to train the model.  
  
abstract static load (fname: str)  
    Load an optimizer from disk and return it :param fname: the filename where the optimizer is serialized  
    :return: The loaded optimizer  
  
abstract save (fname: str) → None  
    Save the optimizer to a file :param fname - the filename to save the optimizer to
```

---

```
abstract test(model: torch.nn.Module, clean_test_data: torch.utils.data.Dataset, triggered_test_data: torch.utils.data.Dataset, clean_test_triggered_labels_data: torch.utils.data.Dataset, torch_dataloader_kwarg) → dict
```

Perform whatever tests desired on the model with clean data and triggered data, return a dictionary of results. :param model: (torch.nn.Module) Trained Pytorch model :param clean\_test\_data: (CSVDataset) Object containing clean test data :param triggered\_test\_data: (CSVDataset or None) Object containing triggered test data, None if triggered data

was not provided for testing

#### Parameters

- **clean\_test\_triggered\_labels\_data** – triggered part of the training dataset but with correct labels; see DataManger.load\_data for more information.
- **torch\_dataloader\_kwarg** – additional arguments to pass to PyTorch’s DataLoader class

#### Returns

(dict) Dictionary of test accuracy results. Required key, value pairs are:

clean\_accuracy: (float in [0, 1]) classification accuracy on clean data clean\_n\_total: (int) number of examples in clean test set

#### The following keys are optional, but should be used if triggered test data was provided

triggered\_accuracy: (float in [0, 1]) classification accuracy on triggered data triggered\_n\_total: (int) number of examples in triggered test set

NOTE: This list may be augmented in the future to allow for additional test data collection.

```
abstract train(model: torch.nn.Module, data: torch.utils.data.Dataset, progress_bar_disable: bool, torch_dataloader_kwarg: dict = None) -> (torch.nn.Module, typing.Sequence[trojai.modelgen.training_statistics.EpochStatistics], <class 'int'>)
```

Train the given model using parameters in self.training\_params :param model: (torch.nn.Module) The untrained Pytorch model :param data: (CSVDataset) Object containing training data, output 0 from TrojaiDataManager.load\_data() :param progress\_bar\_disable: (bool) Don’t display the progress bar if True :param torch\_dataloader\_kwarg: additional arguments to pass to PyTorch’s DataLoader class :return: (torch.nn.Module, EpochStatistics) trained model, a sequence of EpochStatistics objects (one for

each epoch), and the # of epochs with which the model was trained (useful for early stopping).

## trojai.modelgen.runner module

```
class trojai.modelgen.runner.Runner(runner_cfg: trojai.modelgen.config.RunnerConfig, persist_metadata: dict = None)
```

Bases: object

Fundamental unit of model generation, which trains a model as specified in a RunnerConfig object.

**run()** → None

Trains a model and saves it and the associated model statistics

```
trojai.modelgen.runner.add_numerical_extension(path, filename)
```

```
trojai.modelgen.runner.try_force_json(x)
```

Tries to make a value JSON serializable

```
trojai.modelgen.runner.try_serialize(d, u)
```

**trojai.modelgen.torchtext\_optimizer module**

```
class trojai.modelgen.torchtext_optimizer.TorchTextOptimizer(optimizer_cfg: tro-
                                                               jai.modelgen.config.TorchTextOptimizerCon-
                                                               = None)
Bases: trojai.modelgen.optimizer_interface.OptimizerInterface

An optimizer for training and testing LSTM models. Currently in a prototype state.

convert_dataset_to_dataiterator(dataset: trojai.modelgen.datasets.CSVTextDataset,
                                 batch_size: int = None) → torchtext.data.iterator.Iterator

get_cfg_as_dict() → dict
    Return a dictionary with key/value pairs that describe the parameters used to train the model.

get_device_type() → str
    Returns a string representing the device used to train the model

static load(fname: str) → trojai.modelgen.optimizer_interface.OptimizerInterface
    Reconstructs an TorchTextOptimizer, by loading the configuration used to construct the original Torch-
    TextOptimizer, and then creating a new TorchTextOptimizer object from the saved configuration :param
    fname: The filename of the saved TorchTextOptimizer :return: an TorchTextOptimizer object

save(fname: str) → None
    Saves the configuration object used to construct the TorchTextOptimizer. NOTE: because the TorchTex-
    tOptimizer object itself is not persisted, but rather the
        TorchTextOptimizerConfig object, the state of the object does not persist!

    Parameters fname – the filename to save the TorchTextOptimizer's configuration.

test(model: torch.nn.Module, clean_data: trojai.modelgen.datasets.CSVTextDataset, trig-
      gered_data: trojai.modelgen.datasets.CSVTextDataset, clean_test_triggered_labels_data:
      trojai.modelgen.datasets.CSVTextDataset, progress_bar_disable: bool = False,
      torch_dataloader_kwargs: dict = None) → dict
    Test the trained network :param model: the trained module to run the test data through :param
    clean_data: the clean Dataset :param triggered_data: the triggered Dataset, if None, not computed :param
    clean_test_triggered_labels_data: triggered part of the training dataset but with correct labels; see
        DataManger.load_data for more information.

    Parameters
        • progress_bar_disable – if True, disables the progress bar
        • torch_dataloader_kwargs – additional arguments to pass to PyTorch's Dat-
            aLoader class

    Returns a dictionary of the statistics on the clean and triggered data (if applicable)

train(net: torch.nn.Module, dataset: trojai.modelgen.datasets.CSVTextDataset, progress_bar_disable:
      bool = False, torch_dataloader_kwargs: dict = None) -> (torch.nn.Module, typ-
      ing.Sequence[trojai.modelgen.training_statistics.EpochStatistics], <class 'int'>)
    Train the network. :param net: the model to train :param dataset: the dataset to train the network on :param
    progress_bar_disable: if True, disables the progress bar :param torch_dataloader_kwargs: additional argu-
    ments to pass to PyTorch's DataLoader class :return: the trained network, list of EpochStatistics objects,
    and the # of epochs on which teh net was trained
```

---

**train\_epoch**(*model*: *torch.nn.Module*, *train\_loader*: *torchtext.data.iterator.Iterator*, *val\_loader*: *torchtext.data.iterator.Iterator*, *epoch\_num*: *int*, *progress\_bar\_disable*: *bool* = *False*)  
Runs one epoch of training on the specified model

**Parameters**

- **model** – the model to train for one epoch
- **train\_loader** – a DataLoader object pointing to the training dataset
- **val\_loader** – a DataLoader object pointing to the validation dataset
- **epoch\_num** – the epoch number that is being trained
- **progress\_bar\_disable** – if True, disables the progress bar

**Returns** a list of statistics for batches where statistics were computed

```
static train_val_dataset_split(dataset: torchtext.data.Dataset, split_amt: float,
                               val_data_transform: Callable, val_label_transform:
                               Callable) -> (torchtext.data.Dataset, torchtext.data.Dataset)
```

Splits a torchtext dataset (of type: *torchtext.data.Dataset*) into train/test. NOTE: although this has the same functionality as *default\_optimizer.train\_val\_dataset\_split*, it works with a *torchtext.data.Dataset* object rather than *torch.utils.data.Dataset*.

**TODO:** [ ] - specify random seed to torch splitter**Parameters**

- **dataset** – the dataset to be split
- **split\_amt** – fraction specifying the validation dataset size relative to the whole. 1-split\_amt will be the size of the training dataset
- **val\_data\_transform** – (function: any -> any) how to transform the validation data to fit into the desired model and objective function
- **val\_label\_transform** – (function: any -> any) how to transform the validation labels

**Returns** a tuple of the train and validation datasets**trojai.modelgen.training\_statistics module**

```
class trojai.modelgen.training_statistics.BatchStatistics(batch_num: int,
                                                          batch_train_accuracy:
                                                          float, batch_train_loss:
                                                          float)
```

Bases: *object*

Represents the statistics collected from training a batch NOTE: this is currently unused!

```
get_batch_num()
get_batch_train_acc()
get_batch_train_loss()
set_batch_train_acc(acc)
set_batch_train_loss(loss)
```

```
class trojai.modelgen.training_statistics.EpochStatistics(epoch_num,      train-
                                                        ing_stats=None,    val-
                                                        idation_stats=None,
                                                        batch_training_stats=None)
Bases: object
Contains the statistics computed for an Epoch

add_batch(batches:           Union[trojai.modelgen.training_statistics.BatchStatistics,      Se-
                                         sequence[trojai.modelgen.training_statistics.BatchStatistics]])
```

```
get_batch_stats()
```

```
get_epoch_num()
```

```
get_epoch_training_stats()
```

```
get_epoch_validation_stats()
```

```
validate()
```

```
class trojai.modelgen.training_statistics.EpochTrainStatistics(train_acc: float,
                                                               train_loss:      float)
Bases: object
Defines the training statistics for one epoch of training

get_train_acc()
```

```
get_train_loss()
```

```
validate()
```

```
class trojai.modelgen.training_statistics.EpochValidationStatistics(val_clean_acc,
                                                                     val_clean_loss,
                                                                     val_triggered_acc,
                                                                     val_triggered_loss)
Bases: object
Defines the validation statistics for one epoch of training

get_val_acc()
```

```
get_val_clean_acc()
```

```
get_val_clean_loss()
```

```
get_val_loss()
```

```
get_val_triggered_acc()
```

```
get_val_triggered_loss()
```

```
validate()
```

```
class trojai.modelgen.training_statistics.TrainingRunStatistics
Bases: object
Contains the statistics computed for an entire training run, a sequence of epochs TODO:

[ ] - have another function which returns detailed statistics per epoch in an easily serialized manner
```

```
add_best_epoch_val(best_epoch)
```

```
add_epoch(epoch_stats:           Union[trojai.modelgen.training_statistics.EpochStatistics,      Se-
                                         sequence[trojai.modelgen.training_statistics.EpochStatistics]])
```

---

```
add_num_epochs_trained(num_epochs)
autopopulate_final_summary_stats()
```

**Uses the information from the final epoch's final batch to auto-populate the following statistics:**  
final\_train\_acc final\_train\_loss final\_val\_acc final\_val\_loss

```
get_epochs_stats()
```

```
get_summary()
```

Returns a dictionary of the summary statistics from the training run

```
save_detailed_stats_to_disk(fname: str) → None
```

Saves all batch statistics for every epoch as a CSV file

**Parameters** `fname` – filename to save the detailed information to

**Returns** None

```
save_summary_to_json(json_fname: str) → None
```

Saves the training summary to a JSON file

```
set_final_clean_data_n_total(n)
```

```
set_final_clean_data_test_acc(acc)
```

```
set_final_clean_data_triggered_label_n(n)
```

```
set_final_clean_data_triggered_label_test_acc(acc)
```

```
set_final_train_acc(acc)
```

```
set_final_train_loss(loss)
```

```
set_final_triggered_data_n_total(n)
```

```
set_final_triggered_data_test_acc(acc)
```

```
set_final_val_clean_acc(acc)
```

```
set_final_val_clean_loss(loss)
```

```
set_final_val_combined_acc(acc)
```

```
set_final_val_combined_loss(loss)
```

```
set_final_val_triggered_acc(acc)
```

```
set_final_val_triggered_loss(loss)
```

```
trojai.modelgen.training_statistics.logger = <Logger trojai.modelgen.training_statistics (W)
```

Contains classes necessary for collecting statistics on the model during training

## trojai.modelgen.uge\_model\_generator module

```
trojai.modelgen.uge_model_generator.ALL_EXEC_PERMISSIONS = 365
```

This file contains all the functionality needed to train models for a Univa Grid Engine (UGE) HPC cluster.

```
class trojai.modelgen.uge_model_generator.UGEModelGenerator(configs:  
    Union[trojai.modelgen.config.ModelGeneratorConfig,  
          Sequence[trojai.modelgen.config.ModelGeneratorConfig],  
          UGEConfig],  
    uge_config: trojai.modelgen.config.UGEConfig,  
    work_directory: str = '/home/docs/uge_model_generator',  
    validate_uge_dirs: bool = True)  
Bases: trojai.modelgen.model_generator_interface.ModelGeneratorInterface
```

Class which generates models utilizing a Univa Grid Engine

**expand\_modelgen\_configs\_to\_process () → Sequence[trojai.modelgen.config.ModelGeneratorConfig]**  
Converts a sequence of ModelGeneratorConfig objects into another sequence of ModelGeneratorConfig objects such that each element in the sequence only creates one model. For example:

Input: cfgs = [cfg1->num\_models=1, cfg2->num\_models=2]. len(cfgs)=2 Output: cfgs = [cfg1->num\_models=1, cfg2->num\_models=1, cfg2->num\_models=1]. len(cfgs)=3

**NOTE: This will lead to multiple configs pointing to the same data on disk. I'm not sure if this is a problem for PyTorch or not, but this is something to investigate if unexpected results arise.**

**Returns** expanded config configuration

**get\_queue\_numjobs\_assignment () → Sequence**

Determine the number of jobs to give to each queue based on UGEConfig :return: a list of tuples, with each tuple containing the queue in index-0, and the number of jobs assigned to that queue in index-1

**run (mock=False) → None**

Run's the actual UGE job. :param mock: if True, then it generates all the necessary scripts but doesn't execute the UGE command :return: None

**validate () → None**

Validate the input configuration

## trojai.modelgen.utils module

**trojai.modelgen.utils.clamp (X, l, u, cuda=True)**

Clamps a tensor to lower bound l and upper bound u. :param X: the tensor to clamp. :param l: lower bound for the clamp. :param u: upper bound for the clamp. :param cuda: whether the tensor should be on the gpu.

**trojai.modelgen.utils.get\_uniform\_delta (shape, eps, requires\_grad=True)**

Generates a torch uniform random matrix of shape within +-eps. :param shape: the tensor shape to create. :param eps: the epsilon bounds 0+-eps for the uniform random tensor. :param requires\_grad: whether the tensor requires a gradient.

**trojai.modelgen.utils.make\_trojai\_model\_dict (model)**

**Create a TrojAI approved dictionary specification of a PyTorch model for saving to a file. E.g. for a trained model**

**‘model’:** save\_dict = make\_trojai\_model\_dict(model) torch.save(save\_dict, filename)

**Parameters** `model` – (torch.nn.Module) The desired model to be saved.

**Returns** (dict) dictionary containing TrojAI approved information about the model, which can also be used for later loading the model.

```
trojai.modelgen.utils.resave_trojai_model_as_dict(file, new_loc=None)
```

**Load a fully serialized Pytorch model (i.e. whole model was saved instead of a specification) and save it as a TrojAI style dictionary specification.**

#### Parameters

- `file` – (str) Location of the file to re-save
- `new_loc` – (str) Where to save the file if replacing the original is not desired

### Module contents

## 6.2 Module contents



## PYTHON MODULE INDEX

### t

trojai, 61  
trojai.datagen, 43  
trojai.datagen.common\_label\_behaviors,  
    23  
trojai.datagen.config, 24  
trojai.datagen.constants, 26  
trojai.datagen.datatype\_xforms, 26  
trojai.datagen.entity, 26  
trojai.datagen.experiment, 26  
trojai.datagen.image\_affine\_xforms, 28  
trojai.datagen.image\_conversion\_utils,  
    29  
trojai.datagen.image\_entity, 30  
trojai.datagen.image\_insert\_utils, 31  
trojai.datagen.image\_size\_xforms, 31  
trojai.datagen.image\_triggers, 33  
trojai.datagen.insert\_merges, 34  
trojai.datagen.instagram\_xforms, 36  
trojai.datagen.label\_behavior, 37  
trojai.datagen.merge\_interface, 37  
trojai.datagen.pipeline, 38  
trojai.datagen.static\_color\_xforms, 38  
trojai.datagen.transform\_interface, 39  
trojai.datagen.utils, 40  
trojai.datagen.xform\_merge\_pipeline, 40  
trojai.modelgen, 61  
trojai.modelgen.architecture\_factory,  
    44  
trojai.modelgen.architectures, 44  
trojai.modelgen.architectures.cifar10\_architectures,  
    43  
trojai.modelgen.architectures.mnist\_architectures,  
    44  
trojai.modelgen.config, 44  
trojai.modelgen.constants, 48  
trojai.modelgen.data\_configuration, 48  
trojai.modelgen.data\_descriptions, 49  
trojai.modelgen.data\_manager, 49  
trojai.modelgen.datasets, 50  
trojai.modelgen.default\_optimizer, 52  
trojai.modelgen.model\_generator, 54

trojai.modelgen.model\_generator\_interface,  
    54  
trojai.modelgen.optimizer\_interface, 54  
trojai.modelgen.runner, 55  
trojai.modelgen.torchtext\_optimizer, 56  
trojai.modelgen.training\_statistics, 57  
trojai.modelgen.uge\_model\_generator, 59  
trojai.modelgen.utils, 60



# INDEX

## A

add\_batch () (*trojai.modelgen.training\_statistics.EpochStatistics*, 58) (in module *trojai.modelgen.training\_statistics*)  
add\_best\_epoch\_val () (tro-  
  *jai.modelgen.training\_statistics.TrainingRunStatistics* (in module *trojai.modelgen.training\_statistics*), 58)  
add\_epoch () (tro-  
  *jai.modelgen.training\_statistics.TrainingRunStatistics* (in module *trojai.modelgen.training\_statistics*), 58)  
add\_num\_epochs\_trained () (tro-  
  *jai.modelgen.training\_statistics.TrainingRunStatistics* (in module *trojai.modelgen.training\_statistics*), 58)  
add\_numerical\_extension () (in module *trojai.modelgen.runner*), 55  
AlexNet (class in *trojai.modelgen.architectures.cifar10\_architectures*), 43  
ALL\_EXEC\_PERMISSIONS (in module *trojai.modelgen.uge\_model\_generator*), 59  
ArchitectureFactory (class in *trojai.modelgen.architecture\_factory*), 44  
autopopulate\_final\_summary\_stats () (tro-  
  *jai.modelgen.training\_statistics.TrainingRunStatistics* (in module *trojai.modelgen.training\_statistics*), 59)

## B

BadNetExample (class in *trojai.modelgen.architectures.mnist\_architectures*), 44  
BatchStatistics (class in *trojai.modelgen.training\_statistics*), 57  
Bottleneck (class in *trojai.modelgen.architectures.cifar10\_architectures*), 43  
build\_vocab () (tro-  
  *jai.modelgen.datasets.CSVTextDataset* (in module *trojai.modelgen.datasets*), 50)

## C

check\_list\_type () (in module *jai.datagen.config*), 25  
check\_non\_negative () (in module *jai.datagen.config*), 25

clamp () (in module *trojai.modelgen.utils*), 60  
ClassicExperiment (class in *trojai.datagen.experiment*), 26  
ConfigInterface (class in *trojai.modelgen.config*), 44  
convert\_dataset\_to\_dataiterator () (tro-  
  *jai.modelgen.torchtext\_optimizer.TorchTextOptimizer* (in module *trojai.modelgen.torchtext\_optimizer*), 56)  
create () (*trojai.datagen.image\_triggers.RandomRectangularPattern* (in module *trojai.datagen.image\_triggers*), 33)  
create () (*trojai.datagen.image\_triggers.RectangularPattern* (in module *trojai.datagen.image\_triggers*), 33)  
create () (*trojai.datagen.image\_triggers.ReverseLambdaPattern* (in module *trojai.datagen.image\_triggers*), 34)  
create\_experiment () (tro-  
  *jai.datagen.experiment.ClassicExperiment* (in module *trojai.datagen.experiment*), 26)  
csv\_dataset\_from\_df () (in module *trojai.modelgen.datasets*), 51  
csv\_textdataset\_from\_df () (in module *trojai.modelgen.datasets*), 51  
CSVDataset (class in *trojai.modelgen.datasets*), 50  
CSVImageDatasetDesc (class in *trojai.modelgen.data\_descriptions*), 49  
CSVTextDataset (class in *trojai.modelgen.datasets*), 50

## D

DataConfiguration (class in *trojai.modelgen.data\_configuration*), 48  
DataDescription (class in *trojai.modelgen.data\_descriptions*), 49  
DataManager (class in *trojai.modelgen.data\_manager*), 49  
DatasetInterface (class in *trojai.modelgen.datasets*), 51  
default\_image\_file\_loader () (in module *trojai.modelgen.datasets*), 51  
DefaultOptimizer (class in *trojai.modelgen.default\_optimizer*), 52

DefaultOptimizerConfig (class in tro- do () (trojai.datagen.insert\_merges.InsertRandomLocationNonzeroAlpha  
jai.modelgen.config), 44 method), 35  
DefaultSoftToHardFn (class in tro- do () (trojai.datagen.insert\_merges.InsertRandomWithMask  
jai.modelgen.config), 45 method), 35  
DenseNet (class in tro- do () (trojai.datagen.insert\_merges.RandomInsertTextMerge  
jai.modelgen.architectures.cifar10\_architectures), method), 35  
43 do () (trojai.datagen.instagram\_xforms.FilterXForm  
DenseNet121 () (in module tro- method), 36  
jai.modelgen.architectures.cifar10\_architectures), do () (trojai.datagen.label\_behavior.LabelBehavior  
43 method), 37  
DenseNet161 () (in module tro- do () (trojai.datagen.merge\_interface.ImageMerge  
jai.modelgen.architectures.cifar10\_architectures), method), 37  
43 do () (trojai.datagen.merge\_interface.Merge method),  
DenseNet169 () (in module tro- 38  
jai.modelgen.architectures.cifar10\_architectures), do () (trojai.datagen.merge\_interface.TextMerge  
43 method), 38  
DenseNet201 () (in module tro- do () (trojai.datagen.static\_color\_xforms.GrayscaleToRGBXForm  
jai.modelgen.architectures.cifar10\_architectures), method), 38  
43 do () (trojai.datagen.static\_color\_xforms.RGBAtoRGB  
densenet\_cifar () (in module tro- method), 39  
jai.modelgen.architectures.cifar10\_architectures), do () (trojai.datagen.static\_color\_xforms.RGBtoRGBA  
43 method), 39  
do () (trojai.datagen.common\_label\_behaviors.StaticTarget do () (trojai.datagen.transform\_interface.ImageTransform  
method), 23 method), 39  
do () (trojai.datagen.common\_label\_behaviors.WrappedAdd do () (trojai.datagen.transform\_interface.TextTransform  
method), 23 method), 39  
do () (trojai.datagen.datatype\_xforms.ToTensorXForm do () (trojai.datagen.transform\_interface.Transform  
method), 26 method), 39  
do () (trojai.datagen.image\_affine\_xforms.PerspectiveXForm  
method), 28 E  
do () (trojai.datagen.image\_affine\_xforms.RandomPerspectiveXForm do () (trojai.datagen.transform\_interface.ImageTransform  
method), 28 method), 39  
do () (trojai.datagen.image\_affine\_xforms.RandomRotateXForm EntityMappingConfig (class in tro-  
method), 28 jai.modelgen.config), 45  
do () (trojai.datagen.image\_affine\_xforms.RotateXForm Entity (class in trojai.datagen.entity), 26  
method), 28 EpochStatistics (class in tro-  
do () (trojai.datagen.image\_affine\_xforms.UniformScaleXForm jai.modelgen.training\_statistics), 57  
method), 28 EpochTrainStatistics (class in tro-  
do () (trojai.datagen.image\_size\_xforms.Pad method), jai.modelgen.training\_statistics), 58  
31 EpochValidationStatistics (class in tro-  
do () (trojai.datagen.image\_size\_xforms.RandomPadToSize expand\_modelgen\_configs\_to\_process()  
method), 31 (trojai.modelgen.uge\_model\_generator.UGEModelGenerator  
do () (trojai.datagen.image\_size\_xforms.RandomResize method), 60  
method), 32 F  
do () (trojai.datagen.image\_size\_xforms.RandomSubCrop filter () (trojai.datagen.instagram\_xforms.FilterXForm  
method), 32 method), 36  
do () (trojai.datagen.image\_size\_xforms.Resize filter () (trojai.datagen.instagram\_xforms.GothamFilterXForm  
method), 32 method), 36  
do () (trojai.datagen.insert\_merges.FixedInsertTextMerge filter () (trojai.datagen.instagram\_xforms.KelvinFilterXForm  
method), 34 method), 36  
do () (trojai.datagen.insert\_merges.InsertAtLocation filter () (trojai.datagen.instagram\_xforms.LomoFilterXForm  
method), 34 method), 36  
do () (trojai.datagen.insert\_merges.InsertAtRandomLocation filter () (trojai.datagen.instagram\_xforms.NashvilleFilterXForm  
method), 34 method), 37

filter() (trojai.datagen.instagram\_xforms.NoOpFilterXForm.get\_data() (trojai.datagen.image\_triggers.RectangularPattern method), 37  
filter() (trojai.datagen.instagram\_xforms.ToasterXForm.get\_data() (trojai.datagen.image\_triggers.ReverseLambdaPattern method), 37  
FilterXForm (class in tro- get\_data\_description() (tro-  
jai.datagen.instagram\_xforms), 36 jai.modelgen.datasets.CSVDataset method),  
FixedInsertTextMerge (class in tro- 50  
jai.datagen.insert\_merges), 34 get\_data\_description() (tro-  
forward() (trojai.modelgen.architectures.cifar10\_architectures.AlexNet.get\_data\_description() (tro-  
method), 43 jai.modelgen.datasets.CSVTextDataset method), 50  
forward() (trojai.modelgen.architectures.cifar10\_architectures.BattleNet.get\_data\_description() (tro-  
method), 43 jai.modelgen.datasets.DatasetInterface  
forward() (trojai.modelgen.architectures.cifar10\_architectures.DenseNet.get\_device\_type() (tro-  
method), 43  
forward() (trojai.modelgen.architectures.cifar10\_architectures.TrajOptimizer.get\_default\_optimizer\_config() (tro-  
method), 43 method), 44  
forward() (trojai.modelgen.architectures.mnist\_architectures.BadNetExample.get\_device\_type() (tro-  
method), 44 jai.modelgen.config.OptimizerConfigInterface  
forward() (trojai.modelgen.architectures.mnist\_architectures.ModernDenseNet.get\_device\_type() (tro-  
method), 44 jai.modelgen.config.TorchTextOptimizerConfig  
method), 46  
**G**  
GenericImageEntity (class in tro- get\_device\_type() (tro-  
jai.datagen.image\_entity), 30 jai.modelgen.default\_optimizer.DefaultOptimizer  
get\_batch\_num() (tro- method), 52  
jai.modelgen.training\_statistics.BatchStatistics  
get\_batch\_stats() (tro- get\_device\_type() (tro-  
jai.modelgen.training\_statistics.EpochStatistics jai.modelgen.optimizer\_interface.OptimizerInterface  
method), 54  
method), 58  
get\_batch\_train\_acc() (tro- get\_device\_type() (tro-  
jai.modelgen.training\_statistics.BatchStatistics jai.modelgen.torchtext\_optimizer.TorchTextOptimizer  
method), 57 method), 56  
get\_batch\_train\_loss() (tro- get\_epoch\_num() (tro-  
jai.modelgen.training\_statistics.BatchStatistics jai.modelgen.training\_statistics.EpochStatistics  
method), 57 method), 58  
get\_cfg\_as\_dict() (tro- get\_epoch\_training\_stats() (tro-  
jai.modelgen.config.TrainingConfig method), 58  
method), 47 jai.modelgen.training\_statistics.EpochStatistics  
get\_cfg\_as\_dict() (tro- get\_epoch\_validation\_stats() (tro-  
jai.modelgen.default\_optimizer.DefaultOptimizer method), 58  
method), 52 jai.modelgen.training\_statistics.EpochStatistics  
get\_cfg\_as\_dict() (tro- get\_epochs\_stats() (tro-  
jai.modelgen.optimizer\_interface.OptimizerInterface get\_mask() (trojai.datagen.image\_entity.GenericImageEntity  
method), 59 method), 30  
method), 54  
get\_cfg\_as\_dict() (tro- get\_mask() (trojai.datagen.image\_entity.ImageEntity  
jai.modelgen.torchtext\_optimizer.TorchTextOptimizer method), 30  
method), 56 get\_mask() (trojai.datagen.image\_triggers.RandomRectangularPattern  
method), 33  
get\_data() (trojai.datagen.entity.Entity method), 26  
get\_data() (trojai.datagen.image\_entity.GenericImageEntity.get\_mask() (trojai.datagen.image\_triggers.RectangularPattern  
method), 30 method), 33  
get\_data() (trojai.datagen.image\_triggers.RandomRectangularPattern.get\_mask() (trojai.datagen.image\_triggers.ReverseLambdaPattern  
method), 33 method), 34  
get\_predefined\_perspective\_xform\_matrix()

(in module *jai.datagen.image\_affine\_xforms*), 29  
get\_queue\_numjobs\_assignment() (tro- InsertAtLocation (class *jai.datagen.insert\_merges*), 34  
method), 60 (tro- InsertAtRandomLocation (class *jai.datagen.insert\_merges*), 34  
InsertRandomLocationNonzeroAlpha (class in tro-  
get\_summary() (tro- tro- *jai.modelgen.training\_statistics.TrainingRunStatistics*.InsertRandomWithMask (class in tro-  
method), 59 (tro- *jai.datagen.insert\_merges*), 35  
get\_train\_acc() (tro-  
method), 58 (tro- KelvinFilterXForm (class in tro-  
get\_train\_loss() (tro- *jai.datagen.instagram\_xforms*), 36  
method), 58 (tro- L  
get\_uniform\_delta() (in module tro- LabelBehavior (class in tro-  
*jai.modelgen.utils*), 60 (tro- *jai.datagen.label\_behavior*), 37  
get\_val\_acc() (tro- load() (trojai.modelgen.config.DefaultOptimizerConfig  
method), 58 static method), 44  
get\_val\_clean\_acc() (tro- load() (trojai.modelgen.config.ModelGeneratorConfig  
method), 58 static method), 45  
get\_val\_clean\_loss() (tro- load() (trojai.modelgen.config.OptimizerConfigInterface  
method), 58 static method), 45  
get\_val\_loss() (tro- load() (trojai.modelgen.default\_optimizer.DefaultOptimizer  
method), 58 static method), 52  
get\_val\_triggered\_acc() (tro- load() (trojai.modelgen.torchtext\_optimizer.TorchTextOptimizer  
method), 58 static method), 56  
get\_val\_triggered\_loss() (tro- load\_data() (trojai.modelgen.data\_manager.DataManager  
method), 49  
GothamFilterXForm (class in tro- loggers (in module tro-  
*jai.datagen.instagram\_xforms*), 36 tro- *jai.datagen.common\_label\_behaviors*), 23  
gray\_to\_rgb() (in module tro- logger (in module trojai.datagen.config), 25  
*jai.datagen.image\_conversion\_utils*), 29 tro- logger (in module trojai.datagen.datatype\_xforms), 26  
GrayscaleToRGBXForm (class in tro- logger (in module trojai.datagen.entity), 26  
*jai.datagen.static\_color\_xforms*), 38 tro- logger (in module trojai.datagen.experiment), 27  
logger (in module trojai.datagen.image\_affine\_xforms), 29  
logger (in module trojai.datagen.image\_conversion\_utils), 29  
logger (in module trojai.datagen.image\_entity), 30  
logger (in module trojai.datagen.image\_size\_xforms), 32  
logger (in module trojai.datagen.image\_triggers), 34  
logger (in module trojai.datagen.insert\_merges), 35  
logger (in module trojai.datagen.static\_color\_xforms), 39  
logger (in module trojai.datagen.utils), 40  
logger (in module trojai.datagen.xform\_merge\_pipeline), 40  
logger (in module trojai.modelgen.config), 48  
|  
identity\_function() (in module tro-  
*jai.modelgen.config*), 48 tro-  
identity\_transform() (in module tro-  
*jai.modelgen.datasets*), 51 tro-  
ImageDataConfiguration (class in tro-  
*jai.modelgen.data\_configuration*), 48 tro-  
ImageEntity (class in trojai.datagen.image\_entity), 30 tro-  
ImageMerge (class in trojai.datagen.merge\_interface), 37 tro-  
ImageTransform (class in tro-  
*jai.datagen.transform\_interface*), 39 tro-

```
logger (in module trojan.datagen.pipeline), 38
      jai.modelgen.data_configuration), 49
process () (trojan.datagen.pipeline.Pipeline method),
logger (in module trojan.modelgen.datasets), 52
process () 38
logger (in module trojan.modelgen.training_statistics),
process () (trojan.datagen.xform_merge_pipeline.XFormMerge
      59
method), 40
LomoFilterXForm (class in trojan.datagen.instagram_xforms), 36
process_xform_list () (in module trojan.datagen.utils), 40
```

M

make\_trojai\_model\_dict() (in module trojai.modelgen.utils), 60

Merge (class in trojai.datagen.merge\_interface), 38

ModdedLeNet5Net (class in trojai.modelgen.architectures.mnist\_architectures), 44

modelgen\_cfg\_to\_runner\_cfg() (in module trojai.modelgen.config), 48

ModelGenerator (class in trojai.modelgen.model\_generator), 54

ModelGeneratorConfig (class in trojai.modelgen.config), 45

ModelGeneratorInterface (class in trojai.modelgen.model\_generator\_interface), 54

modify\_clean\_image\_dataset() (in module trojai.datagen.xform\_merge\_pipeline), 40

modify\_clean\_text\_dataset() (in module trojai.datagen.xform\_merge\_pipeline), 41

N

```
NashvilleFilterXForm      (class      in      tro-
                           jai.datagen.instagram_xforms), 36
new_architecture()      (tro-
                           jai.modelgen.architecture_factory.ArchitectureF
                           method), 44
NoOpFilterXForm         (class      in      tro-
                           jai.datagen.instagram_xforms), 37
normalization_from_rgb() (in      module      tro-
                           jai.datagen.image_conversion_utils), 29
normalization_to_rgb()  (in      module      tro-
                           jai.datagen.image_conversion_utils), 29
```

0

OptimizerConfigInterface (class in tro-  
jai.modelgen.config), 45  
OptimizerInterface (class in tro-  
jai.modelgen.optimizer\_interface), 54

P

Pipeline (*class in trojai.datagen.pipeline*), 38  
process () (*trojai.datagen.pipeline.Pipeline method*),  
          38  
process () (*trojai.datagen.xform\_merge\_pipeline.XFormMerge*  
              *method*), 40  
process\_xform\_list ()    (*in module tro-*  
                          *jai.datagen.utils*), 40

R

RANDOM\_STATE\_DRAW\_LIMIT (in module trojai.datagen.constants), 26

RandomInsertTextMerge (class in trojai.datagen.insert\_merges), 35

RandomPadToSize (class in trojai.datagen.image\_size\_xforms), 31

RandomPerspectiveXForm (class in trojai.datagen.image\_affine\_xforms), 28

RandomRectangularPattern (class in trojai.datagen.image\_triggers), 33

RandomResize (class in trojai.datagen.image\_size\_xforms), 32

RandomRotateXForm (class in trojai.datagen.image\_affine\_xforms), 28

RandomSubCrop (class in trojai.datagen.image\_size\_xforms), 32

RectangularPattern (class in trojai.datagen.image\_triggers), 33

ReportingConfig (class in trojai.modelgen.config), 46

resave\_trojai\_model\_as\_dict () (in module trojai.modelgen.utils), 61

Resize (class in trojai.datagen.image\_size\_xforms), 32

ReverseLambdaPattern (class in trojectory jai.datagen.image\_triggers), 33

rgb\_to\_rgba () (in module trojai.datagen.image\_conversion\_utils), 30

rgba\_to\_rgb () (in module trojai.datagen.image\_conversion\_utils), 30

RGBAtoRGB (class in trojai.datagen.static\_color\_xforms), 38

RGBtoRGBA (class in trojai.datagen.static\_color\_xforms), 39

RotateXForm (class in trojai.datagen.image\_affine\_xforms), 28

run () (trojai.modelgen.model\_generator.ModelGenerator method), 54

run () (trojai.modelgen.model\_generator\_interface.ModelGenerator method), 54

run () (trojai.modelgen.runner.Runner method), 55

run () (trojai.modelgen.uge\_model\_generator.UGEModelGenerator method), 60

Runner (class in trojai.modelgen.runner), 55

RunnerConfig (class in trojai.modelgen.config), 46

S

save () (trojai.modelgen.config.DefaultOptimizerConfig  
method), 44

save () (trojai.modelgen.config.ModelGeneratorConfig  
method), 45

save () (trojai.modelgen.config.OptimizerConfigInterface  
method), 45

save () (trojai.modelgen.config.TorchTextOptimizerConfig  
method), 47

save () (trojai.modelgen.default\_optimizer.DefaultOptimizer  
method), 52

save () (trojai.modelgen.optimizer\_interface.OptimizerInterface  
method), 54

save () (trojai.modelgen.torchtext\_optimizer.TorchTextOptimizer  
method), 56

save\_detailed\_stats\_to\_disk () (tro-  
jai.modelgen.training\_statistics.TrainingRunStatistics  
method), 59

save\_summary\_to\_json () (tro-  
jai.modelgen.training\_statistics.TrainingRunStatistics  
method), 59

set\_batch\_train\_acc () (tro-  
jai.modelgen.training\_statistics.BatchStatistics  
method), 57

set\_batch\_train\_loss () (tro-  
jai.modelgen.training\_statistics.BatchStatistics  
method), 57

set\_data\_description () (tro-  
jai.modelgen.datasets.CSVDataset  
method), 50

set\_data\_description () (tro-  
jai.modelgen.datasets.CSVTextDataset  
method), 50

set\_data\_description () (tro-  
jai.modelgen.datasets.DatasetInterface  
method), 51

set\_embedding\_vectors\_cfg () (tro-  
jai.modelgen.data\_configuration.TextDataConfiguration  
method), 48

set\_final\_clean\_data\_n\_total () (tro-  
jai.modelgen.training\_statistics.TrainingRunStatistics  
method), 59

set\_final\_clean\_data\_test () (tro-  
jai.modelgen.training\_statistics.TrainingRunStatistics  
method), 59

set\_final\_clean\_data\_triggered\_label\_n () (tro-  
jai.modelgen.training\_statistics.TrainingRunStatistics  
method), 59

set\_final\_clean\_data\_triggered\_label\_test () (tro-  
jai.modelgen.training\_statistics.TrainingRunStatistics  
method), 59

set\_final\_train\_acc () (tro-  
jai.modelgen.training\_statistics.TrainingRunStatistics  
method), 59

set\_final\_train\_loss () (tro-  
jai.modelgen.training\_statistics.TrainingRunStatistics  
method), 59

set\_final\_triggered\_data\_n\_total () (tro-  
jai.modelgen.training\_statistics.TrainingRunStatistics  
method), 59

set\_final\_triggered\_data\_test\_acc () (tro-  
jai.modelgen.training\_statistics.TrainingRunStatistics  
method), 59

set\_final\_val\_clean\_acc () (tro-  
jai.modelgen.training\_statistics.TrainingRunStatistics  
method), 59

set\_final\_val\_clean\_loss () (tro-  
jai.modelgen.training\_statistics.TrainingRunStatistics  
method), 59

set\_final\_val\_combined\_acc () (tro-  
jai.modelgen.training\_statistics.TrainingRunStatistics  
method), 59

set\_final\_val\_combined\_loss () (tro-  
jai.modelgen.training\_statistics.TrainingRunStatistics  
method), 59

set\_final\_val\_triggered\_acc () (tro-  
jai.modelgen.training\_statistics.TrainingRunStatistics  
method), 59

set\_final\_val\_triggered\_loss () (tro-  
jai.modelgen.training\_statistics.TrainingRunStatistics  
method), 59

setup\_optimizer\_generator () (tro-  
jai.modelgen.config.RunnerConfig  
method), 46

sort\_key () (trojai.modelgen.datasets.CSVTextDataset  
static method), 50

split\_val\_clean\_trig () (in module tro-  
jai.modelgen.default\_optimizer), 53

StaticTarget (class in tro-  
jai.datagen.common\_label\_behaviors), 23

subset\_clean\_df\_by\_labels () (in module tro-  
jai.datagen.xform\_merge\_pipeline), 42

T

test () (trojai.modelgen.default\_optimizer.DefaultOptimizer  
method), 52

test () (trojai.modelgen.optimizer\_interface.OptimizerInterface  
method), 54

test () (trojai.modelgen.torchtext\_optimizer.TorchTextOptimizer  
method), 56

TextDataConfiguration (class in tro-  
jai.modelgen.data\_configuration), 48

TextMerge (class in trojai.datagen.merge\_interface), 38

TextTransform (class in tro-  
jai.datagen.transform\_interface), 39

MasterXForm (class in tro-  
jai.datagen.instagram\_xforms), 37

TorchTextOptimizer (class in *trojai.modelgen.torchtext\_optimizer*), 56  
TorchTextOptimizerConfig (class in *trojai.modelgen.config*), 46  
ToTensorXForm (class in *trojai.datagen.datatype\_xforms*), 26  
train() (*trojai.modelgen.default\_optimizer.DefaultOptimizer*.*method*), 52  
train() (*trojai.modelgen.optimizer\_interface.OptimizerInterface*.*method*), 55  
train() (*trojai.modelgen.torchtext\_optimizer.TorchTextOptimizer* (*module*)), 40  
train\_epoch() (*trojai.modelgen.default\_optimizer.DefaultOptimizer*.*method*), 52  
train\_epoch() (*trojai.modelgen.torchtext\_optimizer.TorchTextOptimizer*.*method*), 56  
train\_val\_dataset\_split() (*in module trojai.modelgen.default\_optimizer*), 53  
train\_val\_dataset\_split() (*trojai.modelgen.torchtext\_optimizer.TorchTextOptimizer*.*static method*), 57  
TrainingConfig (class in *trojai.modelgen.config*), 47  
TrainingRunStatistics (class in *trojai.modelgen.training\_statistics*), 58  
Transform (class in *trojai.datagen.transform\_interface*), 39  
Transition (class in *trojai.modelgen.architectures.cifar10\_architectures*), 43  
trojai (*module*), 61  
trojai.datagen (*module*), 43  
trojai.datagen.common\_label\_behaviors (*module*), 23  
trojai.datagen.config (*module*), 24  
trojai.datagen.constants (*module*), 26  
trojai.datagen.datatype\_xforms (*module*), 26  
trojai.datagen.entity (*module*), 26  
trojai.datagen.experiment (*module*), 26  
trojai.datagen.image\_affine\_xforms (*module*), 28  
trojai.datagen.image\_conversion\_utils (*module*), 29  
trojai.datagen.image\_entity (*module*), 30  
trojai.datagen.image\_insert\_utils (*module*), 31  
trojai.datagen.image\_size\_xforms (*module*), 31  
trojai.datagen.image\_triggers (*module*), 33  
trojai.datagen.insert\_merges (*module*), 34  
trojai.datagen.instagram\_xforms (*module*), 36  
trojai.datagen.label\_behavior (*module*), 37  
trojai.datagen.merge\_interface (*module*), 37  
trojai.datagen.pipeline (*module*), 38  
trojai.datagen.static\_color\_xforms (*module*), 38  
trojai.datagen.transform\_interface (*module*), 39  
trojai.datagen.utils (*module*), 40  
trojai.datagen.xform\_merge\_pipeline  
trojai.modelgen (*module*), 61  
trojai.modelgen.architecture\_factory (*module*), 44  
trojai.modelgen.architectures (*module*), 44  
trojai.modelgen.architectures.cifar10\_architectures (*module*), 43  
trojai.modelgen.architectures.mnist\_architectures (*module*), 44  
trojai.modelgen.config (*module*), 44  
trojai.modelgen.constants (*module*), 48  
trojai.modelgen.data\_configuration (*module*), 48  
trojai.modelgen.data\_descriptions (*module*), 49  
trojai.modelgen.data\_manager (*module*), 49  
trojai.modelgen.datasets (*module*), 50  
trojai.modelgen.default\_optimizer (*module*), 52  
trojai.modelgen.model\_generator (*module*), 54  
trojai.modelgen.model\_generator\_interface (*module*), 54  
trojai.modelgen.optimizer\_interface (*module*), 54  
trojai.modelgen.runner (*module*), 55  
trojai.modelgen.torchtext\_optimizer (*module*), 56  
trojai.modelgen.training\_statistics (*module*), 57  
trojai.modelgen.uge\_model\_generator (*module*), 59  
trojai.modelgen.utils (*module*), 60  
TrojAIConfig (class in *trojai.datagen.config*), 24  
try\_force\_json() (*in module trojai.modelgen.runner*), 55  
try\_serialize() (*in module trojai.modelgen.runner*), 55

**U**

UGEConfig (class in *trojai.modelgen.config*), 47  
UGEModelGenerator (class in *trojai.modelgen.uge\_model\_generator*), 59

UGEQueueConfig (*class in trojai.modelgen.config*), 47  
UniformScaleXForm (*class in trojai.datagen.image\_affine\_xforms*), 28

**V**

VALID\_DEVICES (*in module jai.modelgen.constants*), 48  
valid\_locations () (*in module jai.datagen.image\_insert\_utils*), 31  
VALID\_LOSS\_FUNCTIONS (*in module jai.modelgen.constants*), 48  
VALID\_OPTIMIZERS (*in module jai.modelgen.constants*), 48  
validate () (*trojai.datagen.config.TrojAIConfig* method), 24  
validate () (*trojai.datagen.config.ValidInsertLocationsConfig* method), 24  
validate () (*trojai.datagen.config.XFormMergePipelineConfig* method), 25  
validate () (*trojai.datagen.config.TrojAIConfig* method), 24  
validate () (*trojai.modelgen.config.EarlyStoppingConfig* method), 45  
validate () (*trojai.modelgen.config.ModelGeneratorConfig* method), 45  
validate () (*trojai.modelgen.config.ReportingConfig* method), 46  
validate () (*trojai.modelgen.config.RunnerConfig* method), 46  
validate () (*trojai.modelgen.config.TorchTextOptimizerConfig* method), 47  
validate () (*trojai.modelgen.config.TrainingConfig* method), 47  
validate () (*trojai.modelgen.config.UGEConfig* method), 47  
validate () (*trojai.modelgen.config.UGEQueueConfig* method), 48  
validate () (*trojai.modelgen.data\_configuration.TextDataConfiguration* method), 48  
validate () (*trojai.modelgen.data\_manager.DataManager* method), 50  
validate () (*trojai.modelgen.model\_generator.ModelGenerator* method), 54  
validate () (*trojai.modelgen.training\_statistics.EpochStatistics* method), 58  
validate () (*trojai.modelgen.training\_statistics.EpochTrainStatistics* method), 58  
validate () (*trojai.modelgen.training\_statistics.EpochValidationStatistics* method), 58  
validate () (*trojai.modelgen.uge\_model\_generator.UGEModelGenerator* method), 60  
validate\_model\_generator\_interface\_input ()  
    (*in module jai.modelgen.model\_generator\_interface*), 54

**W**

validate\_optimizer () (*jai.modelgen.config.RunnerConfig* static method), 46  
validate\_regenerate\_mode () (*jai.datagen.config.XFormMergePipelineConfig* method), 25  
ValidInsertLocationsConfig (*class in trojai.datagen.config*), 24

**X**

WrappedAdd (*class in jai.datagen.common\_label\_behaviors*), 23

**XFormMerge** (*class in jai.datagen.xform\_merge\_pipeline*), 40

**XFormMergePipelineConfig** (*class in trojai.datagen.config*), 24