

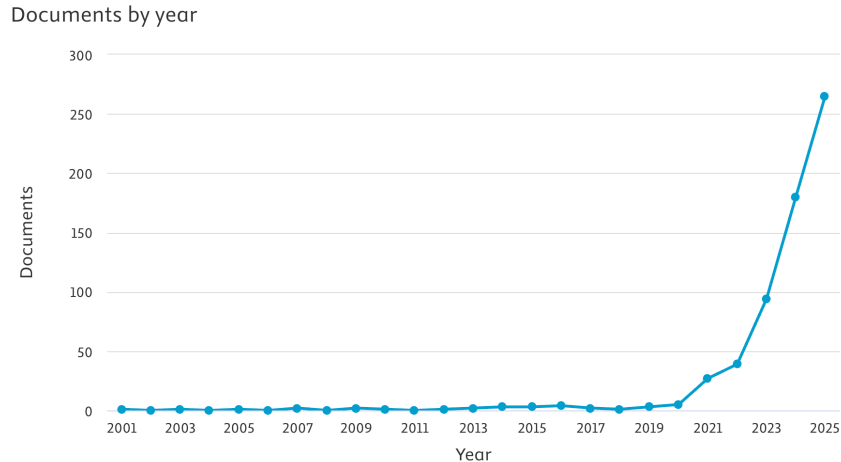


Forgetting Right: Towards Reproducible Benchmarks in Machine Unlearning with ERASURE

Andrea D'Angelo

Università degli Studi dell'Aquila / Italy

Machine Unlearning (MU) is a rapidly growing field, with several publications and methods being published every month.



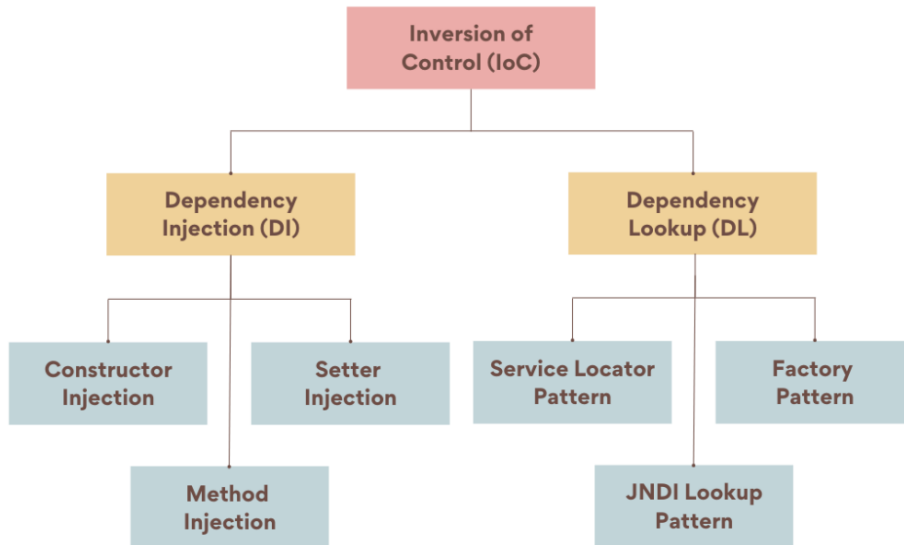
- A quick scopus search shows that the topic is quickly on the rise.

- As often the case for new topics, benchmarking standardization is still lacking.
- We noticed that many works build code only relevant to their specific use case.
- This results in **duplicated efforts** and **low reproducibility**.

- The MU community needs a standardized tool for **reproducible** and **extensible** testing.
- This was our rationale for building **ERASURE**: a modular, extensible framework for Machine Unlearning.

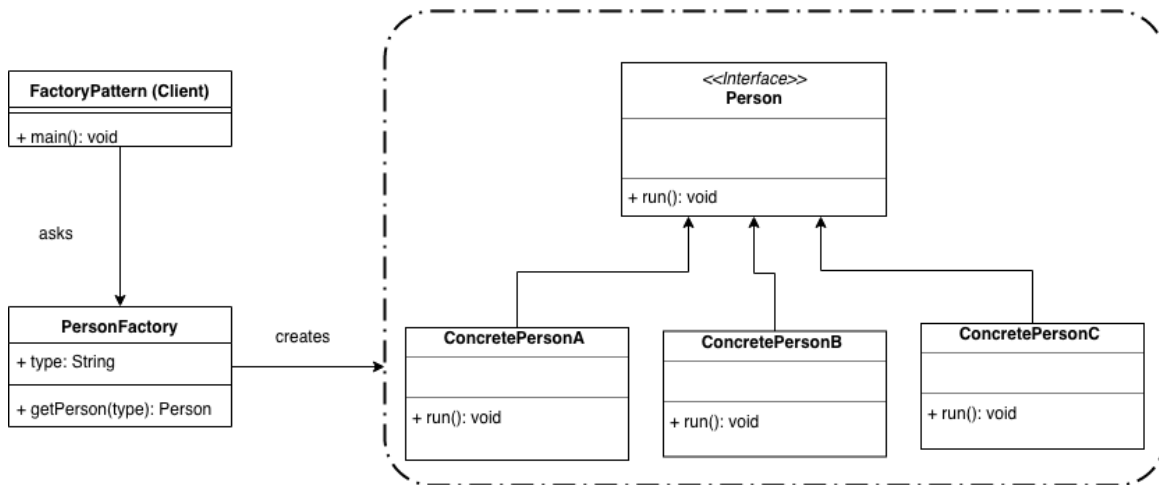


ERASURE is a modular, extensible Machine Unlearning Framework.



- **Inversion of Control (IoC):** shifts control of object creation and execution flow from client code to an external framework. IoC reduces tight coupling by delegating dependency management to an external provider, often through *dependency injection*. ERASURE leverages IoC for dynamic dependency management.

Design Patterns (2)



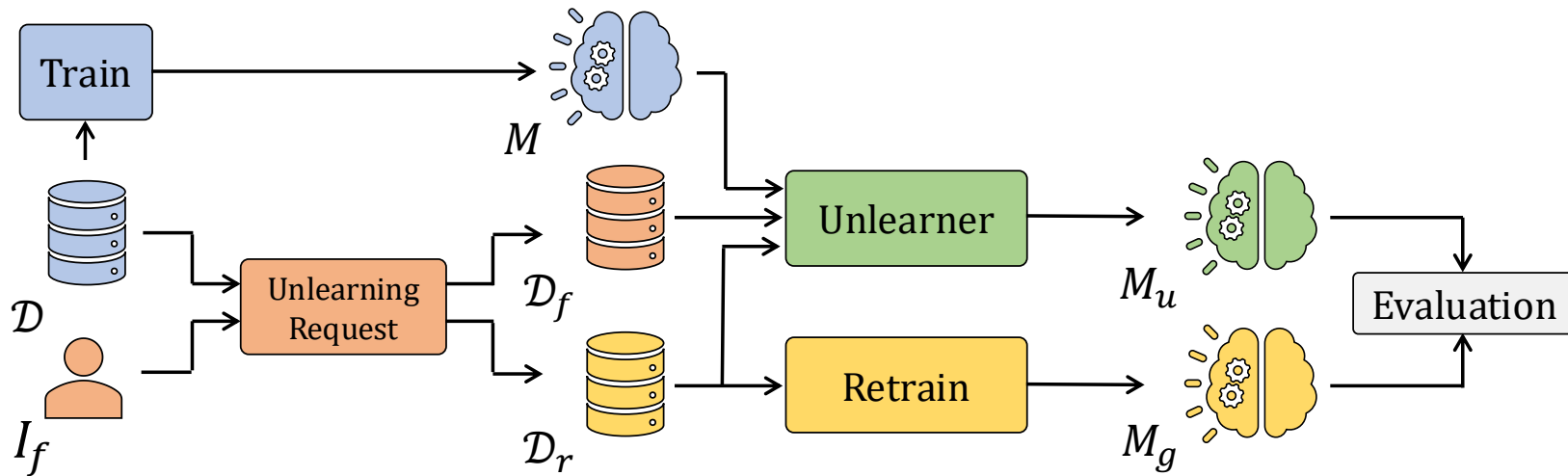
- **Factory Pattern:** abstracts object instantiation. It encapsulates object creation within a dedicated factory class, enabling flexible instantiation. ERASURE uses the Factory Pattern to instantiate Datasets, Models, Unlearners, and everything that is needed for a complete Machine Unlearning workflow.

Machine Unlearning Workflow

7

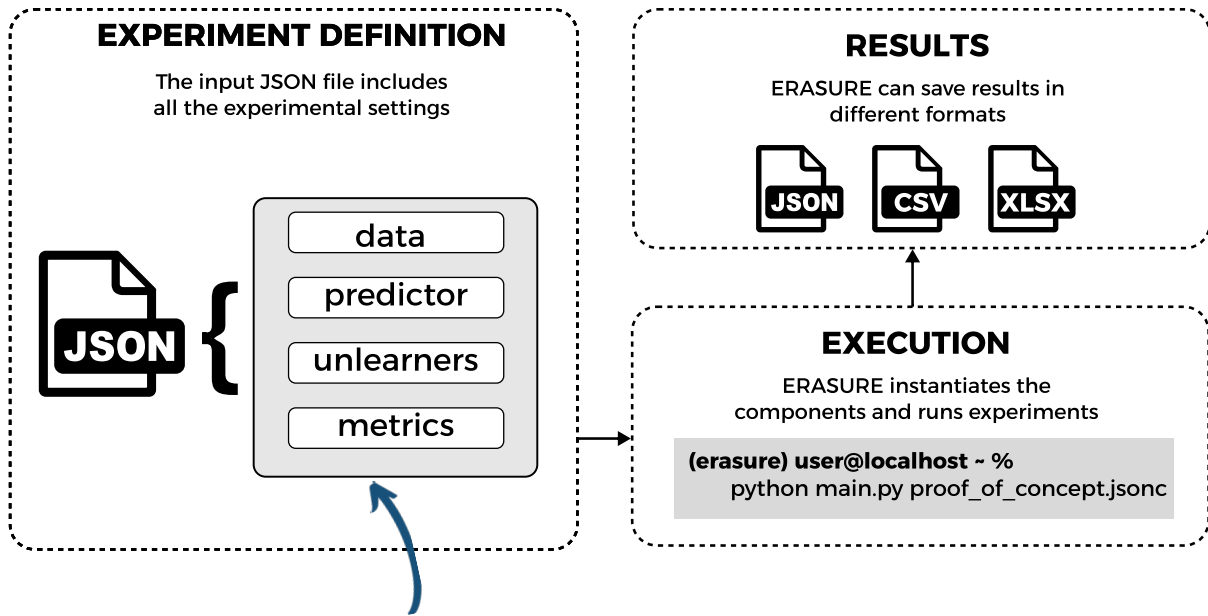
ERASURE was built with the MU workflow at its core.

ERASURE handles all the steps of a typical Machine Unlearning workflow:



ERASURE's Configuration

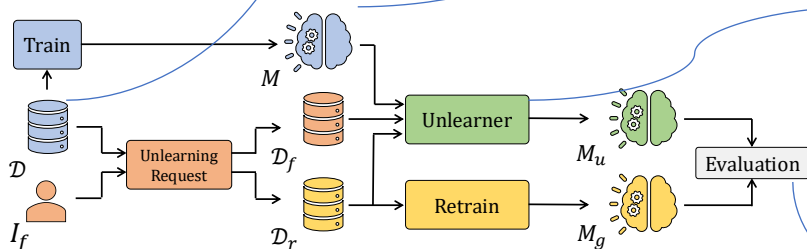
ALL the Experimental settings in ERASURE are defined through a single JSON file, called **Main Configuration**.



The Main Configuration file needs to specify all these components.

Main Configuration

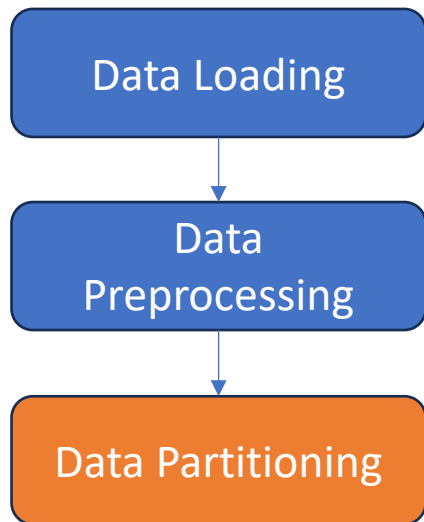
9



```
1 "data": {"class": "erasure.data.<NS>.DatasetManager",
2         "parameters": {
3             "DataSource": { ... },
4             "partitions": [ "p_1", "p_2", ... , "p_n" ]
5         }},
6 "predictor": {"class": "erasure.model.<NS>.TorchModel",
7               "parameters": {
8                 "optimizer": {"class": "torch.optim.Adam"},
9                 "loss_fn": {"class": "torch.nn.CrossEntropyLoss"},
10                "model": {"class": "erasure.<NS>.BERTClassifier"}
11            }},
12 "unlearners": [
13     {"compose_gold": "configs/snippets/u_gold.json"},
14     .
15     .
16     ,
17     {"class": "erasure.<NS>.AdvancedNegGrad",
18       "parameters": {
19         "epochs": 1,
20         "ref_data_retain": "retain",
21         "ref_data_forget": "forget",
22         "optimizer": {"class": "torch.optim.Adam",
23                       "parameters": {"lr": 0.0001}}
24       }
25 ],
26 "evaluator": {
27     "class": "erasure.evaluations.<NS>.Evaluator",
28     "parameters": { "measures": [ ... ] }
29 }
```

These are the same modules of the Machine Unlearning workflow.

Machine Unlearning requires more demanding Data handling with respect to usual downstream tasks.



ERASURE introduces a flexible DataSplitter strategy, enabling datasets to be partitioned by configuration in virtually unlimited ways.

Each DataSplitter defines a specific partitioning logic, such as sampling a fixed percentage of data or selecting samples from a given class.

Since they are executed sequentially, the DataSplitters can reference previously created partitions.

Data Partitioning (2)

11

ERASURE also supports a parameter Z.

Imagine you want to remove all the images of George Clooney from CelebA.

There are the
labels (Y)

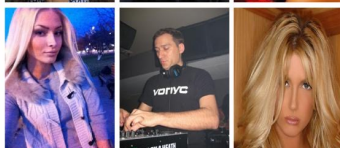
Eyeglasses



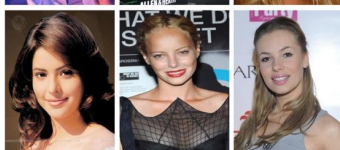
Bangs



Pointy
Nose



Oval Face



Wearing
Hat



Wavy Hair



Mustache



Smiling



There are the
samples (X)

Data Partitioning (3)

12



There is no information on what celebrity is pictured in each photo! If we only pass the plain dataset (X and Y), we will not be able to filter out pictures, which is critical in Machine Unlearning.

Data Partitioning (4)

```
13 "partitions": [  
14   {"class": "erasure.data.<NS>.DataSplitterByZ",  
15     "parameters": {"parts_names": ["forget", "other"],  
16     "z_labels": [1]}},  
17   {"class": "erasure.data.<NS>.DataSplitterPercentage",  
18     "parameters": {"parts_names": ["retain", "test"],  
19     "percentage": 0.8, "ref_data": "other"}},  
20   {"class": "erasure.data.<NS>.DataSplitterConcat",  
    "parameters": {"parts_names": ["train", "-"],  
    "concat_splits": ["retain", "forget"]}}  
],
```

For instance, the first DataSplitter, called **DataSplitterByZ**, splits the data in »forget« and »other« based on the passed z_labels, in this case [1].

Intuitively, if the sample has $z=1$, the sample is put in the forget set.

In **green**, the cascading feature in effect: the second DataSplitter will split only the data from »other«.

```
6 "predictor": {"class": "erasure.model.<NS>.TorchModel",
7  "parameters": {
8    "optimizer": {"class": "torch.optim.Adam"},
9    "loss_fn": {"class": "torch.nn.CrossEntropyLoss"},
10   "model": {"class": "erasure.<NS>.BERTClassifier"}
11  }},
```

ERASURE handles batched training automatically.

All the training parameters are defined in the configuration file, along with the seed for reproducibility.

ERASURE easily handles classes and parameters from external libraries, like Pytorch or sci-kit learn.

The third section of a JSON Main Configuration file is the “unlearners” JSON Array.

Consider, for instance, the Unlearner called **AdvancedNegGrad**.

```
12 "unlearners": [  
13   {"compose_gold" : "configs/snippets/u_gold.json"},  
14   .  
15   .  
16   ,  
17   {"class": "erasure.<NS>.AdvancedNegGrad",  
18     "parameters": {  
19       "epochs": 1,  
20       "ref_data_retain": "retain",  
21       "ref_data_forget": "forget",  
22       "optimizer": {"class": "torch.optim.Adam",  
23                     "parameters": {"lr": 0.0001}}}  
24 ],
```

It requires a series of parameters, like epochs or optimizer, that can be passed directly from the JSON Configuration file.

Implementing new unlearners

To implement a **Custom Unlearner**, one can extend the **TorchUnlearner** class and implement its `__unlearn__(self)` method.

It's that simple.

- `__unlearn__(self)` takes nothing as input but has **access to all the parameters** that were passed through the configuration file, all the dataset and datasets partition, and the original model.
- **The original model is given in isolation**, meaning that all Unlearners has access to their copy. So, the modifications made by one Unlearner will not propagate to other Unlearners.
- The `__unlearn__(self)` method must **return a modified version of the original model**. This returned model will be evaluated lately.

Evaluating Unlearners

ERASURE comes with the most well-known Unlearning metrics out-of-the-box.

Other metrics, like Accuracy, can be instantiated directly from external libraries, like sci-kit learn.

```
25 "evaluator":{  
26   "class": "erasure.evaluations.<NS>.Evaluator",  
27   "parameters": { "measures":[ ... ] } }
```

All metrics have access to both the modified model and the original model, to evaluate differences if they are needed.

ERASURE Out-of-the-box

Datasets

ALL datasets
from
HuggingFace

ALL datasets
from
Torchvision

ALL datasets
from UCI
Repository

ALL datasets
from Torch
Geometric

Unlearners

Gold Model

FineTuning

Advanced
NegGrad

BadTeaching

NegGrad

Eu_k

SRL

FisherForgetting

SalUn

Composite

Scrub

SSD

Unsir

Evaluations

ALL metrics
from sklearn

RunTime

Memory Usage

UMIA

- Extending Erasure is as simple as creating a class with your logic.
- Demo paper at IJCAI for more info:
- CIKM Resource paper coming soon!



Support ERASURE!

20

<https://github.com/aiim-research/ERASURE>



PLEASE SUPPORT US WITH A STAR!

Thank you for your attention

UNIVERSITÀ
DEGLI STUDI
DELL'AQUILA



DISIM
Dipartimento di Ingegneria
e Scienze dell'Informazione
e Matematica