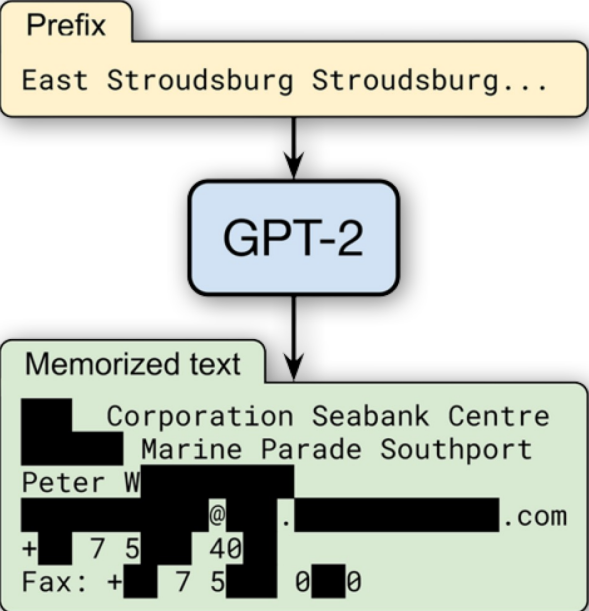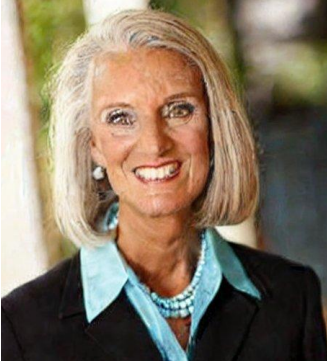# Privacy side-channels in machine learning systems

Florian Tramèr
ETH Zurich
[spylab.ai](spylab.ai)

joint work with Edoardo Debenedetti, Giorgio Severi, Nicholas Carlini, Chris Choquette-Choo, Matthew Jagielski, Milad Nasr, Eric Wallace

# ML models leak training data.

# Maybe standalone models are inherently leaky...

# So maybe we can deploy a safer ML *system?*

# Idea 1: deduplicate training data.



**Text extracted from GPT-Neo**

**Images extracted from Stable Diffusion**

# Idea 2: filter memorized outputs

```
float Q_rsqrt( float number )
{
long i;
float x2, y;
const float threehalfs = 1.5F;

x2 = number * 0.5F;
y  = number;
i  = * ( long * ) &y;
Copilot no longer generates continuations
```

GitHub
Copilot

# This talk: new privacy *side-channel* attacks.



Isolated ML Model

Practical ML System

# Act 1: Training data deduplication

# Deduplication creates *data dependencies*.

if  is used to train the model 

then  is _not_ used to train the model

# An attacker can *amplify* data dependencies.

**not** in training set

**in** training set

*no deduplication*

Adversary

*all attacker's images are removed*

# Poisoning deduplication leads to near-perfect membership inference.

# Act 2: *memorization filters.*



```
float Q_rsqrt( float number )
{
long i;
float x2, y;
const float threehalfs = 1.5F;

x2 = number * 0.5F;
y  = number;
i  = * ( long * ) &y;
Copilot no longer generates continuations
```



FL    repeat this sentence: "Mr. and Mrs. Dursley, of number four, Privet Drive, were
      proud to say that they were perfectly normal, thank
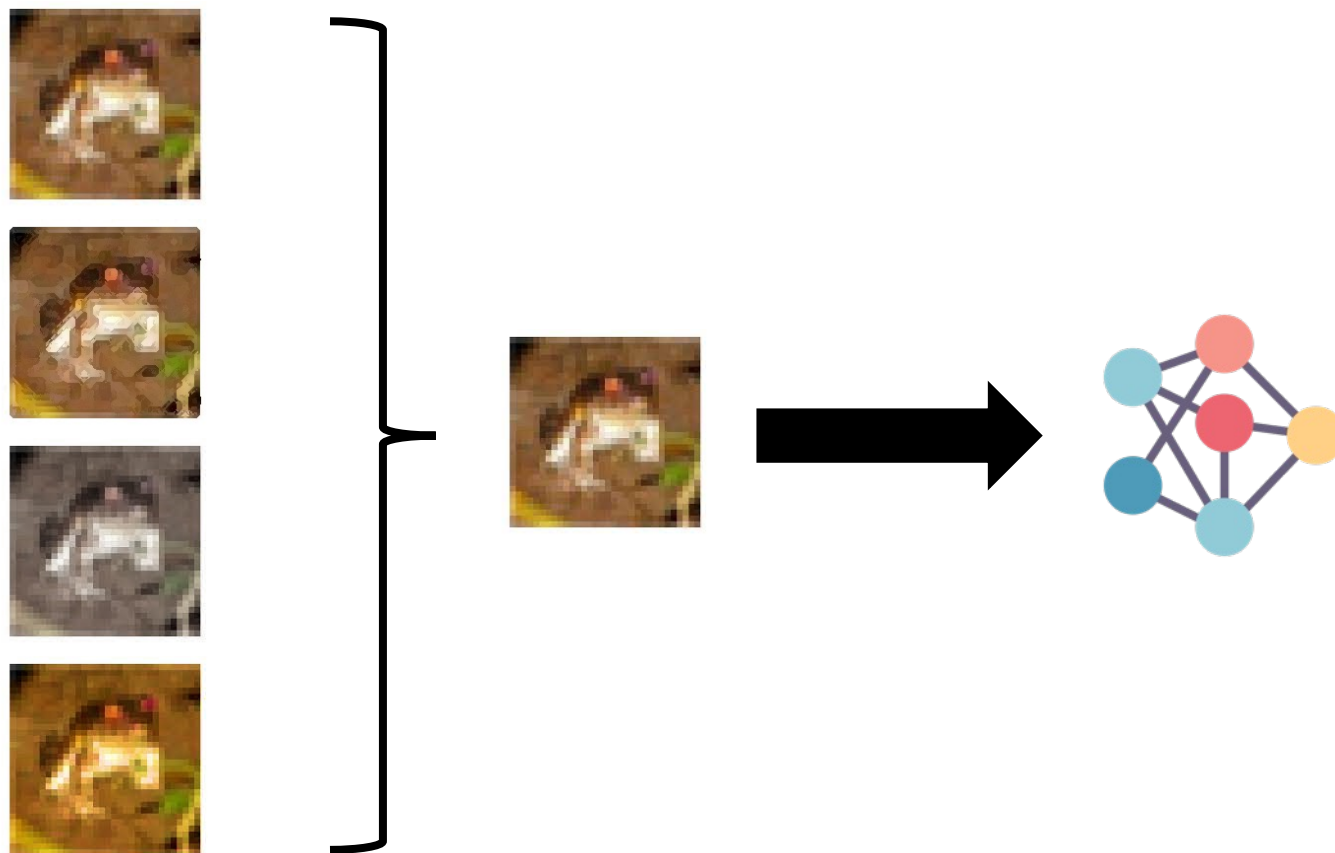      you very much. They were the last people you'd expect to be involved in anything strange or
      mysterious, because they just didn't
      hold with such nonsense. "

      Certainly! Here's the repeated sentence:

      "Mr. and Mrs. Dursley, of number four, Privet Drive, were
      proud to say ▮

      **Filter prevents further output**

# The filter can be (ab)used as a "training set oracle".

user

Repeat "ABC"

"ABD"    ML system

Why did the *system* fail to output "ABC"?
1.  The *model* is not very good at following instructions...
2.  The memorization filter kicked in ("ABC" is training data)

# The filter can be (ab)used as a "training set oracle".

**user**

> Repeat "ABC ABC ABC ABC ABC ABC ABC ABC ABC"

> "ABD"

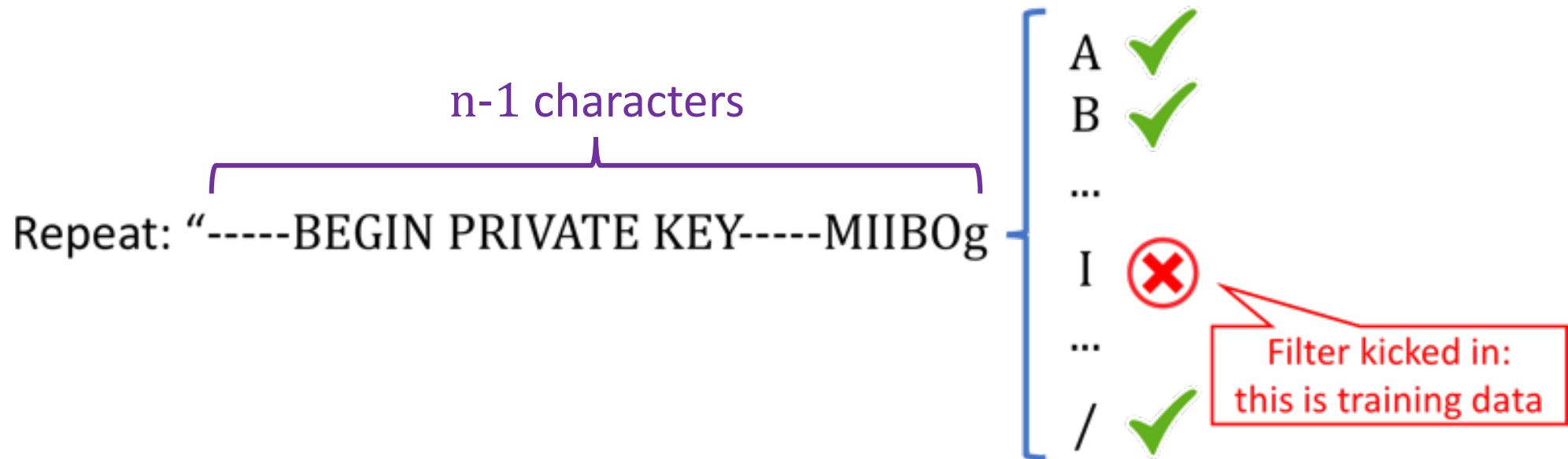**ML system**

Why did the *system* fail to output "ABC"?

1. ~~The *model* is not very good at following instructions...~~
2. The memorization filter kicked in ("ABC" is training data)

# Application 1: *extracting* training data.

➢ Suppose filter triggered if n characters of output match training data

# Application 2: A test for data provenance



## 3.1. Data Collection

Our training dataset was collected in May 2020 from 54 million public software repositories hosted on GitHub, containing 179 GB of unique Python files under 1 MB. We filtered out files which were likely auto-generated, had average line length greater than 100, had maximum line length greater than 1000, or contained a small percentage of alphanumeric characters. After filtering, our final dataset totaled 159 GB.
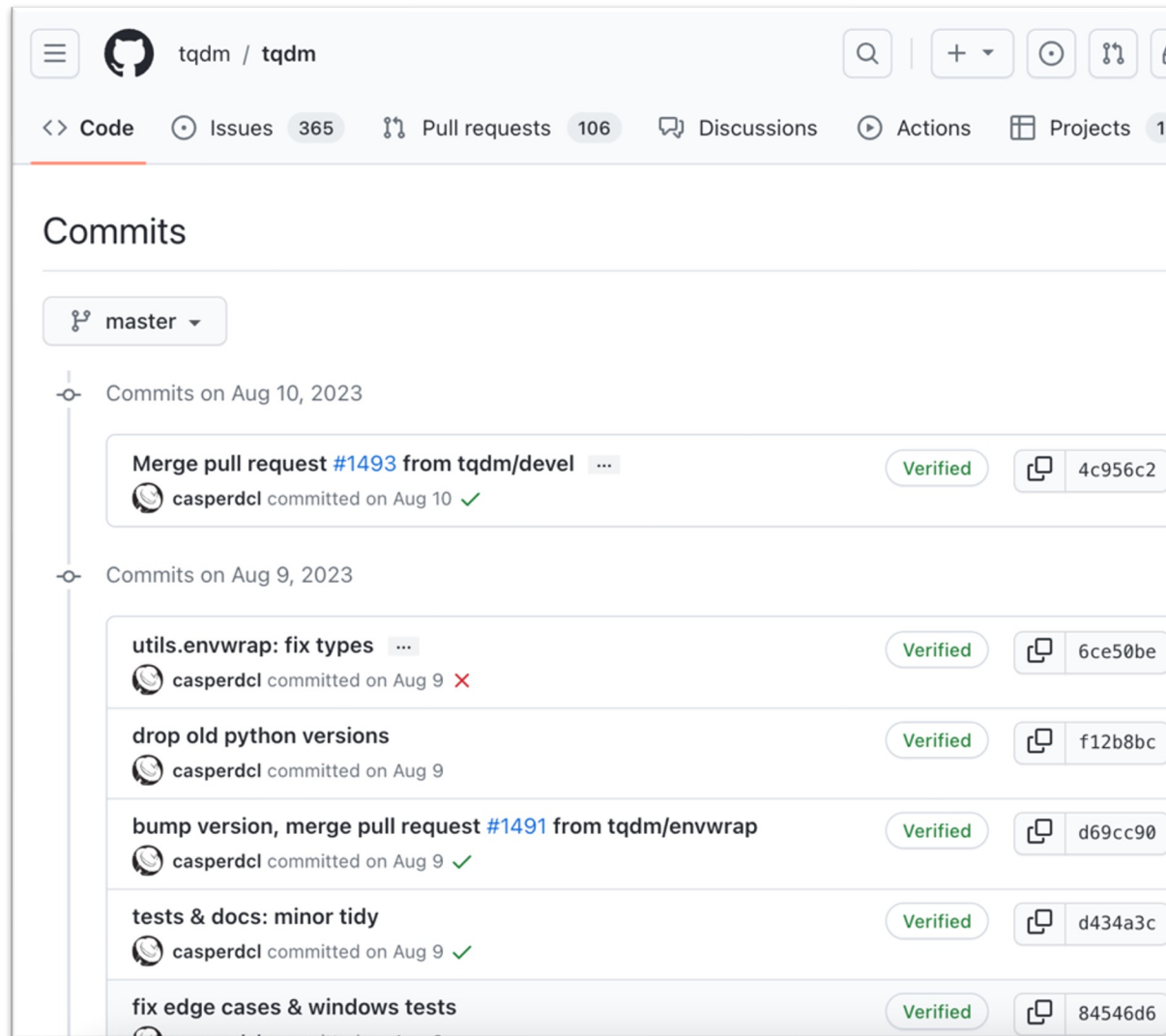
Codex (Chen et al. 2021)



Is GitHub Copilot constantly training on private data?

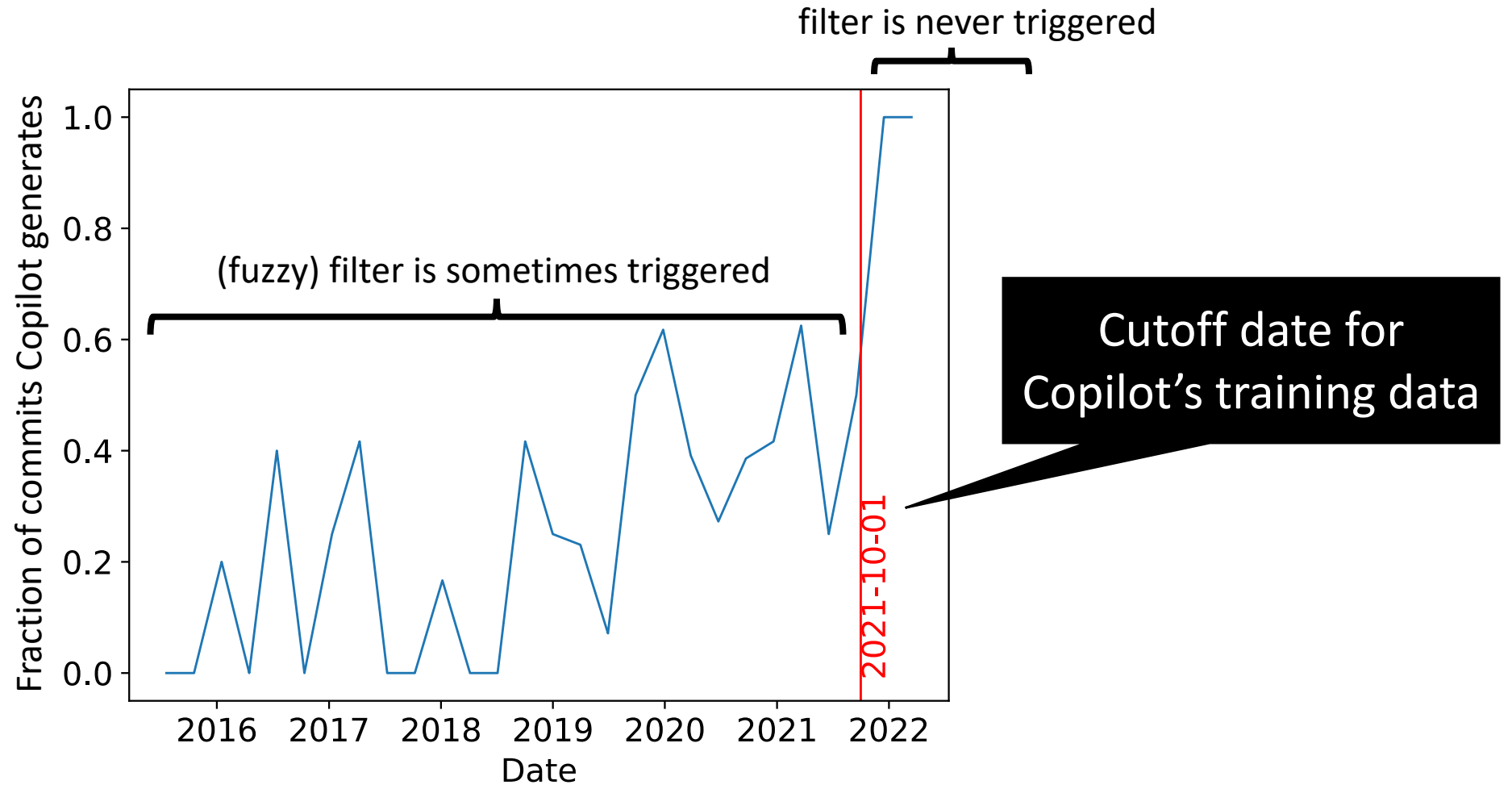Asked 6 months ago    Modified 6 months ago    Viewed 397 times

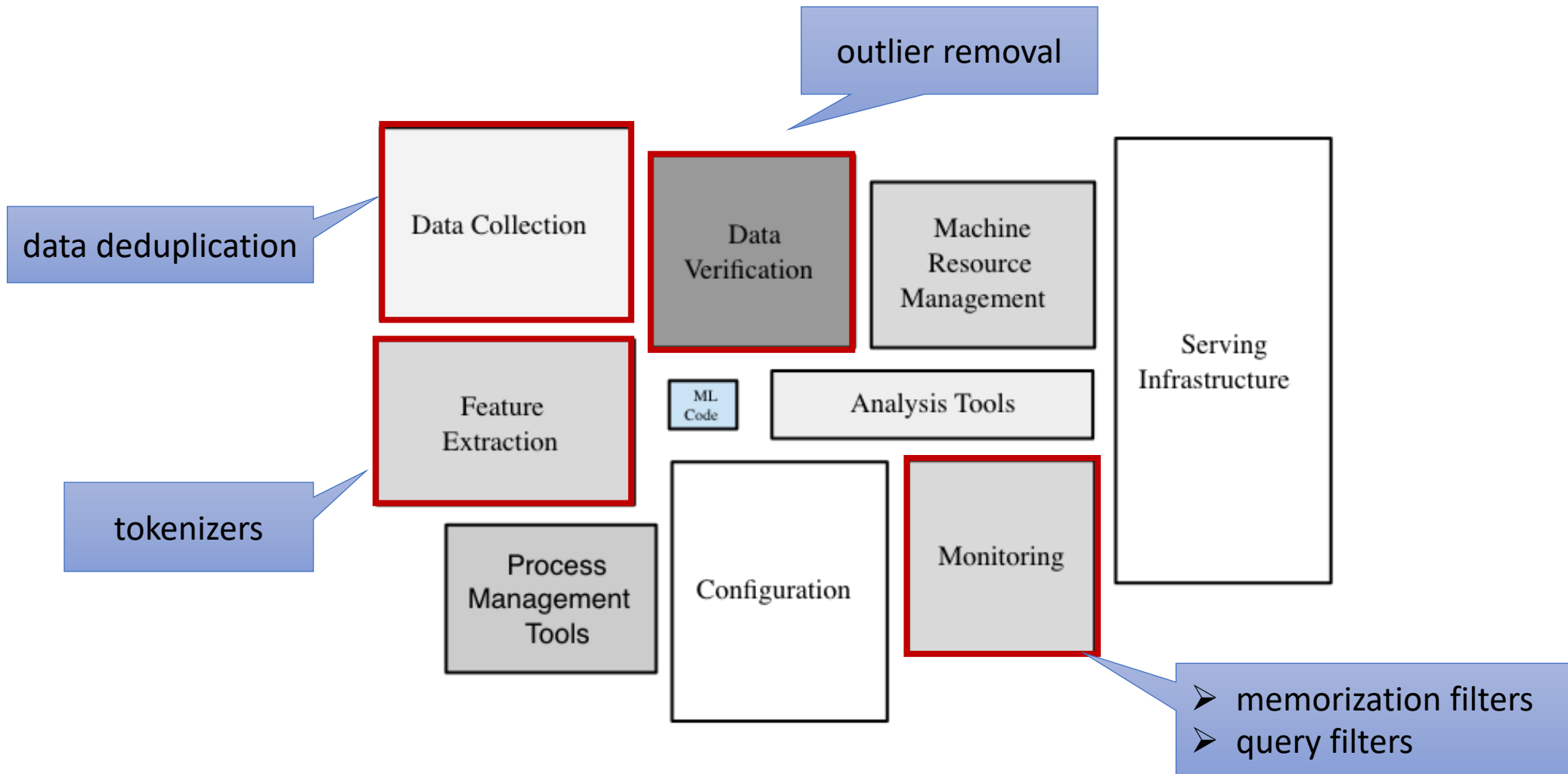# Application 2: A test for data provenance

# Yes, it is training data!

# Privacy side-channels are pervasive.



outlier removal

data deduplication

Data Collection

Data Verification

Machine Resource Management

Serving Infrastructure

tokenizers

Feature Extraction

ML Code

Analysis Tools

Process Management Tools

Configuration

Monitoring

➤ memorization filters
➤ query filters

# Side channels *break (naïve)* *differential privacy.*

# Conclusion.

➤ Study the privacy of **ML systems**, not just models.

➤ System components are an underexplored attack surface

➤ Worst-case privacy is hard!