# Cyber Threat Intelligence
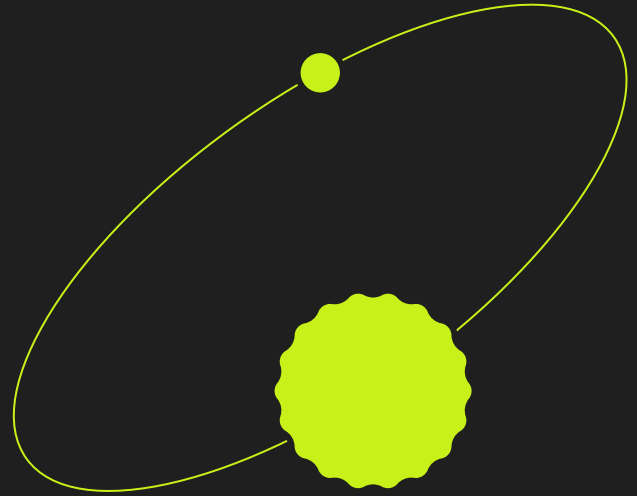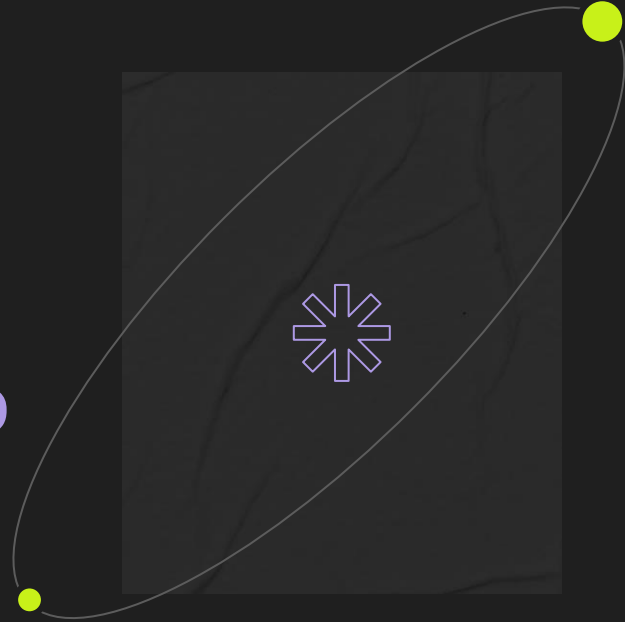
# A.

# What is a data breach crawler?

"Data crawling is a method which involves data mining from different sources"

# How it works?

**Data Breach**

Someone steals personal data.

**Published**

Someone buys the breach and it's published.

**Crawler**

Crawlers store breaches.

# What we use as example:

## Have i been pwned?

Website that allows users to check whether their personal data has been compromised. It collects millions of data; the users can search by entering their phone number or email address.

## Intelligence X

Search engine and data archive. It has unique features: the search works with selectors (email, URLs, IPs, etc.), it searches in places such as the darknet and others. Keeps a historical data archive of results.

# Where did we find them?

**1** **Clear web**

Using search engines and conventional sites.

**2** **Deep web**

Deep web and onion pages, including black-hat sites.

**3** **Undeground forums**

Undeground black hat hacking crime forums.

**4** **Social sites**
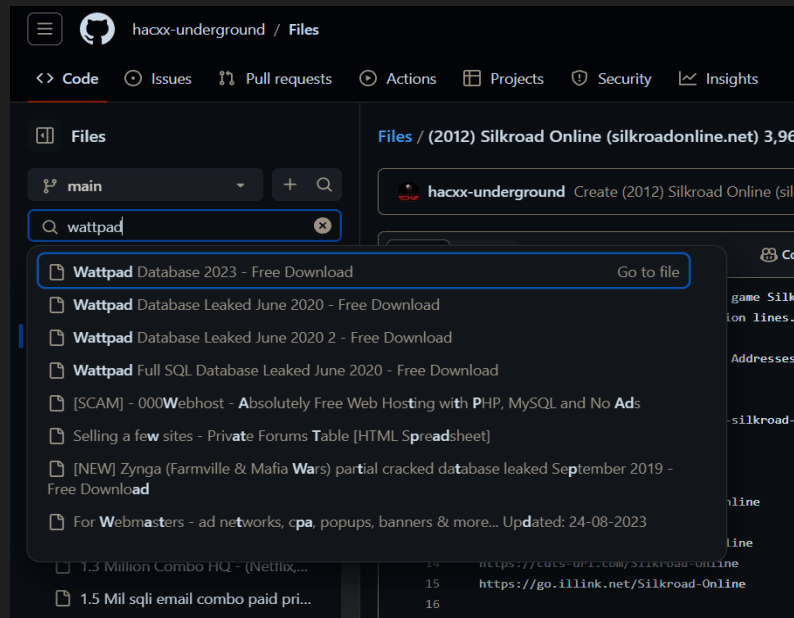
Social network and forums.

# Breaches research

## 1    Clear web

By googling the names of the breached platforms, we can get download links.

They are often embedded in Pastebin pages, Github repositories for researchers, and even in clear web forums post.

# Breaches research

**2**  **Deep web**

On the onion sites of cybercriminal groups can often be found leaked information, often derived from victims who have not paid a ransom.
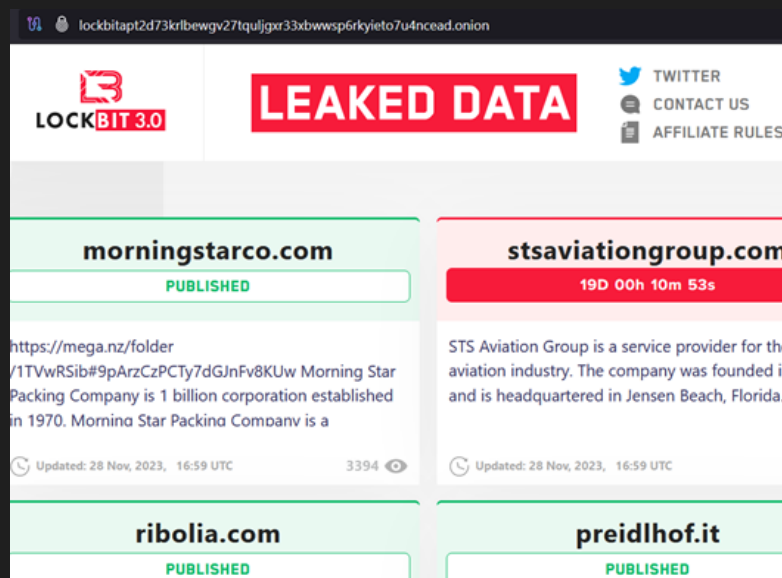
# Breaches research

**3**    **Underground forums**

There are forums where cybercriminals publish breached data, often for free or with a system of tokens. Like Leakbase.io and Breachforums.



threat researchers without breachforums content to screenshot



BreachForums  ›  Leaks  ›  Databases  ›  Forum Announcement

**Forum Announcement: Database Index**

[Owner] pompompurin

March 17, 2022, 05:22 AM

This thread will index all the datasets we have marked as "our CDN.
Please note there are hundreds more unofficial datasets in
This list is not only limited to database breaches, you will fi

If you want to be notified when we add databases to our CI

Sorting options:

- Record Count: [HIGHER-LOWER] [LOWER-HIGHER]
- Date Added: [NEWER-OLDER] [OLDER-NEWER]
- Breach Date: [NEWER-OLDER] [OLDER-NEWER]
- Title: [A-Z] [Z-A]

Bossman

Click here to learn how to get credits.
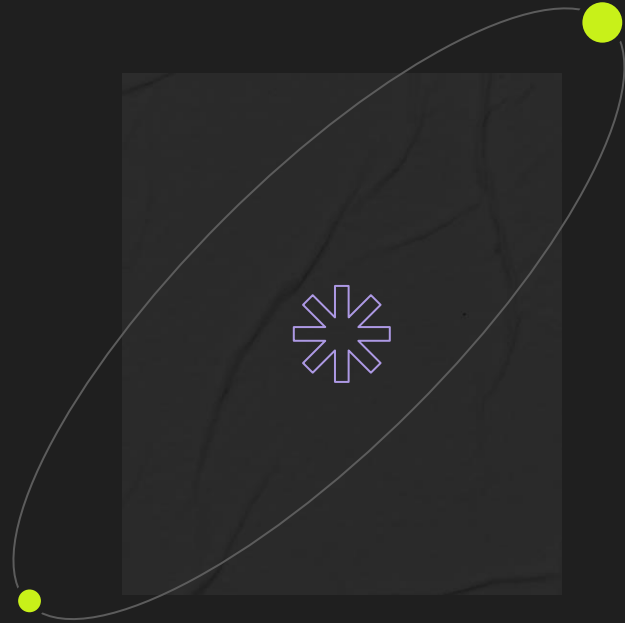basic rules.

# Breaches research

**4**     ## Social networks and forums

Users often post links to leaked data download on places like Reddit, Telegram and Twitter.


r/databreach



**NN Hacking Group** @NNHackingGroup · 5 apr

@emailitsrl BREACHED DA 2 ANNI! Continuate ancora a salvare le password in CHIARO?

- 44 Database
- Dati di 600000 utenti con password in CHIARO!
- SMS/FAX inviati in chiaro
- Messaggi e allegati inviati e ricevuti di tutte le 600k mailbox

nonamemqwl2lslts.onion/emailit.html ❓
#Hacked #0Day



**Facebook Breach | Facebook Leaked Data**
1 254 subscribers

All Country Facebook Breached Data

Facebook Data Leak Files

Yeah Fcuk Facebook 🥸 👾

**VIEW IN TELEGRAM**

# B.

## The project

# Data Breach Selection

Having limited resolutions, the group decided to use specific data breaches. The selection was made based on the information within them (telephone numbers, emails, passwords, etc.) and so that there were no redundancies between the various data breaches.

# Data Breach Selection

## Linkedin

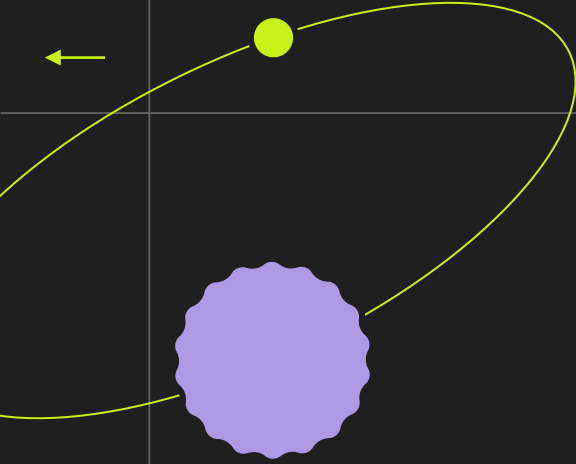ID, Email, password

## Twitter

Email, password

## Stockx

Email, name, surname, address, phone number

## Edmodo

Email, password

## Facebook

ID, name, surname, location, work, email, bithday, phone number

# Functions:

# Project functions

**01** **Password**

Searched by API and our engine.

**02** **Search**

Search for emails and phone numbers within files.

**03** **Identikit**

Creation of the identikit of the chosen person using the information found.

# Password search

We use the "Have I been pwned?" API for the password:

➤ The user enter the password.

➤ We create the HASH of the entered password.

➤ We take the first 5 characters of the HASH and "I have been pwend?" returns a range of HASH.

➤ On this return we select the correct HASH.

➤ The output will show how many times this pass is repeated in the data breaches.

➤ If the password is present in our databases, the table with the users using the entered password will be printed

# Project Versions

Our project has changed for various reasons, here are the versions:

# Version 1

## HTML, CSS, JS

A group was already working in Python, so to differentiate ourselves we decided to use HTML, CSS and JS.

### PROS

Server-less, DBless

Execute on Client

Easy to setup and customized

### CONS

Limited to 1 GB – 10mln record

Depends on client resources

Non-scalable system

**WE ACHIEVED 1 GB / 4 MLN RECORD**

# Version 2

## NODEJS

To use when done so far, we opted to upgrade the system to NodeJS (Server)

## PROS

Horizontal Scalability

Execution Speed I/O

NPM Modules

## CONS

Memory Load Handling

Lack of Native Support for Relational Databases

Asynchronous Complexity Management

**WE ACHIEVED 5 GB / 13 MLN RECORD**

# Version 3

## PYTHON (FLASK + SQLITE)

Noting the difficulties in terms of resources and timing, we decided to base ourselves on flask and sqlite.

| PROS | CONS |
|------|------|
| Rich Ecosystem and Libraries | Performance in certain scenarios |
| Support for Relational Databases | Maintenance and optimization |
| Scalability | Knowledge of different libraries |

**WE ACHIEVED +10 GB / +100 MLN RECORD**

# Version 4

## APACHE SOLR

Apache Solr is the best performing solution for big amount of data, but it requires several hours of study and implementation

<table>
<tr><th>PROS</th><th>CONS</th></tr>
<tr><td>Docker</td><td>Difficult implementation</td></tr>
<tr><td>Multiple Hosts</td><td>Ram usage</td></tr>
<tr><td>Speed</td><td>Data Processing</td></tr>
</table>

**+ 13MLD**

**HARDWARE DEPENDENT**

# Who are you looking for?

write here...

All   E-Mail   Phone Number   Password   First Name   Last Name

# The end

Any questions?