# AIL Project - Practical and Efficient Data-Mining of Chats, Suspicious Websites, Forums and Tor Hidden-Services

CIRCL - Virtual Summer School 2025

⌂ https://ail-project.org/

Alexandre Dulaunoy - alexandre.dulaunoy@circl.lu
Aurelien Thirion - aurelien.thirion@circl.lu
July 17, 2025
CIRCL https://www.circl.lu

# Background

- Over the past years, CIRCL has developed the AIL project[1] to fulfill our needs at CIRCL in intelligence gathering and analysis.
- AIL features an extensible Python-based framework for the **analysis of unstructured information**, collected either through an advanced crawler manager or from various feeders, including social networks and custom feeders.
- The AIL Project is an **open-source** framework[2] comprising various modules designed for the **collection, crawling, digging, and analysis of unstructured data**.
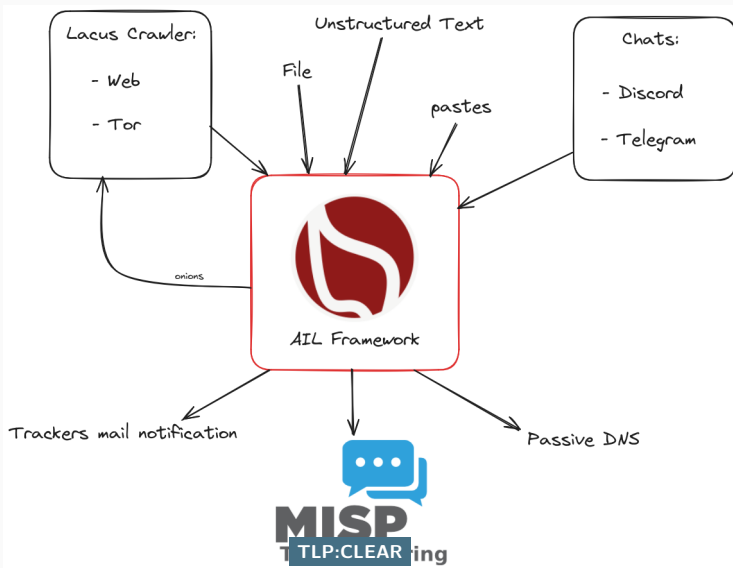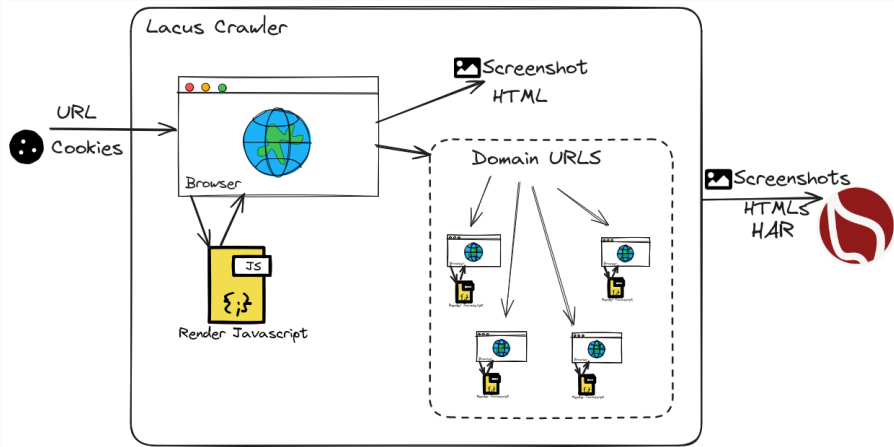


---

[1]https://www.ail-project.org/
[2]https://github.com/ail-project

- Crawling can be a challenging task, for example, gathering all the blog posts from ransomware groups[3], which can be demanding for an analyst.
- AIL offers a crawling feature that can **initiate regular crawls using a standard spawned browser**.



---
[3]https://www.ransomlook.io/

## Use your cookies to login and bypass captcha

### Edit Cookiejar

| Description | Date | UUID | User |
|---|---|---|---|
| 3thxemke2x7hcibu.onion | 2020/03/31 | 90674deb-38fb-4eba-a661-18899ccb3841 | admin@admin.test |

**Edit Description** ✏️  **Add Cookies** ⊕

```
{
    "domain": ".3thxemke2x7hcibu.onion",
    "name": "mybb[lastactive]",
    "path": "/forum/",
    "value": "1583829465"
}
```

```
{
    "domain": ".3thxemke2x7hcibu.onion",
    "name": "loginattempts",
    "path": "/forum/",
    "value": "1"
}
```

```
{
    "domain": ".3thxemke2x7hcibu.onion",
    "name": "sid",
    "path": "/forum/",
    "value": "047ab0cd97ff5bcc77edb6a"
}
```

```
{
    "name": "remember_token",
    "value": "12|58cddd1511d74d341f23"
}
```

```
{
    "domain": ".3thxemke2x7hcibu.onion",
    "name": "mybb[announcements]",
    "path": "/forum/",
    "value": "0"
}
```

# Collection - Automate Collection

- Collecting data from various chat sources can be a **tedious task for analysts**.
- AIL offers a set of feeders (e.g., Telegram, Discord) that can be used to subscribe to chat channels.
- All the **collected messages are then processed and analyzed** within AIL's *processing* and *analysis* stages.

DDoSia Project :



| Name | ID | Created at | First Seen | Last Seen | NB Sub-Channels | Participants |
|------|-----|-----------|-----------|-----------|----------------|-------------|
| DDoSia Project | 2125229770 | 2024-03-07 09:03:02 | 2024-03-07 | 2024-04-12 | 6 | 695 |

Tags: ⊞

🕮 Investigations                                                                              👁 Correlations

## Sub-Channels:

Show 10 ⬍ entries                                                                    Search: [        ]

| Icon | Name | ID | Created at | First Seen | Last Seen | 💬 |
|------|------|-----|-----------|-----------|-----------|-----|
| 🔴 | General | 2125229770/1 | 2024-03-07 06:43:47 | 2024-03-07 | 2024-04-12 | 4121 |
| 🔴 | Поддержка - отвечаем на вопросы<br>Support - answering questions | **TLP:CLEAR** 24-03-07 10:00:24 | 2024-03-07 | 2024-04-12 | 2264 |
| 🔴 | English support | 2125229770/138 | 2024-03-07 08:03:02 | 2024-03-07 | 2024-04-11 | 297 |

- Threat actors are often verbose and frequently **share extensive details in private channels**.

- Many messages contain screenshots and images.

- Text detection and extraction are performed across $80+$ languages using a CRNN (Convolutional Recurrent Neural Network).

- **Enables keyword-based matching and detection**.

## AIL Framework: Features

- Extracting **credit card numbers, credentials, phone numbers, . . .**
- Extracting and validating potential **hostnames**
- Submission to threat sharing and incident response platforms (**MISP** and **FlowIntel**[5])
- **Tagging**[6] for classification and searches
- Terms, sets, regex, and YARA **tracking, occurrences, and history**
- Archives, files, and raw **submission** from the UI
- Correlation engine based on PGP ID, cryptocurrencies, decoded (Base64, . . . ), usernames, cookie names, and many selectors to find relationships
- And many more

---

[5] https://github.com/flowintel/flowintel
[6] Relying on MISP taxonomies and galaxy

## Live Trackers & Retro Hunt

- Search and monitor specific keywords/patterns
  - Automatic tagging
  - Email/webhook notifications
- Track Word
  - ddos
- Track Set
  - booter, ddos, stresser; 2
- Track Regex
  - circl\lu
- **YARA rules**
  - https://github.com/ail-project/ail-yara-rules

# Live Demo

# Dashboard

## Tor and Web Search:

**Type:**

[🧛 Tor] [🌐 Web] [🧛🌐 All]

Tor ▼ | content to Search | [🔍]

## 🧛🌐 Search Domain by name:

| Domain name | [🔍] |

🔵 🔴 Onion Domains
⚪ 🟡 Web Domains

## H Titles Search:

Content Search ▼ | ID or content to Search | [🔍]

🔵 Case Sensitive

## 💬 Chats Search:

**Type:**

[💬 discord] [matrix] [✈ telegram] [💬 All Chats]

discord ▼ | content to Search | [🔍]

## 👤 Usernames Search:

### Sidebar

- 🧛 Tor and Web Search
- 🧛🌐 Search Domain by name
- H Title Search
- 💬 Chats Search
- 👤 Username Search
- ✉ Mail Search
- G GTracking Search
- 🗎 File Name Search

[≡ Toggle Sidebar]

### Navbar

Home | Submit | Tags | Leaks Hunter | Crawlers | Objects | Search | Settings | Log Out

# YARA Tracker

**Certificate**

| | |
|---|---|
| Type | yara |
| Tracked | ail-yara-rules/rules/crypto/certificate.yar |
| Date | 2023/05/12 |
| Level | Global |
| Creator | admin@admin.test |
| First Seen | 2023 / 05 / 12 |
| Last Seen | 2023 / 05 / 31 |
| Tags | |
| Mails | |
| Webhook | |
| Filters | No Filters |
| Objects Match | decoded 0  item 88 |

Edit Tracker ✎  🗑

Yara Rule:

```
rule certificates
{
    meta:
        author = "@KevTheHermit"
        info = "Part of PasteHunter"
        reference = "https://github.com/kevthehermit/PasteHunter"

    strings:
        $ssh_priv = "BEGIN RSA PRIVATE KEY" wide ascii nocase
        $openssh_priv = "BEGIN OPENSSH PRIVATE KEY" wide ascii nocase
        $dsa_priv = "BEGIN DSA PRIVATE KEY" wide ascii nocase
        $ec_priv = "BEGIN EC PRIVATE KEY" wide ascii nocase
        $pgp_priv = "BEGIN PGP PRIVATE KEY" wide ascii nocase
        $pem_cert = "BEGIN CERTIFICATE" wide ascii nocase
        $pkcs7 = "BEGIN PKCS7"

    condition:
        any of them
}
```

📅 2023-05-12     📅 2023-05-31

🔍 Tracked Objects

ail-yara-rules/rules/c...

**TLP:CLEAR**

## test

✓ completed

🔍 Show Objects

| | |
|---|---|
| Date | 2023/05/10 |
| Description | None |
| Tags | |
| Creator | admin@admin.test |
| Filters | `{` `"item": {` `"date_from": "20230304",` `"date_to": "20230601"` `}` `}` |
| Objects Match | item 3 |

```
rule certificates
{
    meta:
        author = "@KevTheHermit"
        info = "Part of PasteHunter"
        reference = "https://github.com/kevthehermit/PasteHunter"

    strings:
        $ssh_priv = "BEGIN RSA PRIVATE KEY" wide ascii nocase
        $openssh_priv = "BEGIN OPENSSH PRIVATE KEY" wide ascii nocase
        $dsa_priv = "BEGIN DSA PRIVATE KEY" wide ascii nocase
        $ec_priv = "BEGIN EC PRIVATE KEY" wide ascii nocase
        $pgp_priv = "BEGIN PGP PRIVATE KEY" wide ascii nocase
        $pem_cert = "BEGIN CERTIFICATE" wide ascii nocase
        $pkcs7 = "BEGIN PKCS7"

    condition:
        any of them
}
```

Show 10 ⌄ entries                                                    Search: [        ]

| Type | Id | Tags |
|---|---|---|
| ● | archive/gist.github.com/2023/04/14/luizmiranda7_3b3d1133a3d3842092c5fc5fb39e84f2.gz | infoleak:automatic-detection="private-key" test23 test12 infoleak:automatic-detection="certificate" |
| ● | submitted/2023/04/20/submitted_cc9190ab-80d2-4d2b-9c9e-97c51e69a855.gz | infoleak:submission="manual" test12 infoleak:automatic-detection="rsa-private-key" infoleak:automatic-detection="vpn-static-key" test23 infoleak:automatic-detection="certificate" infoleak:automatic-detection="onion" |
| ● | archive/gist.github.com/2023/04/13/chipzoller_d8d6d2d737d02ad4fe9d30a897170761.gz | test12 test23 infoleak:automatic-detection="certificate" |

Showing 1 to 3 of 3 entries                                    Previous 1 Next

## Tor Coin Mixer

| | |
|---|---|
| UUID | 9189d0e7c04c47a29f85666e9507e0a5 |
| Creator | admin@admin.test |
| Tags | dark-web:topic="mixer" |
| Date | 2023-05-31 |
| Threat Level | medium |
| Analysis | initial |
| Info | Tor Coin Mixer |
| # Objects | 6 |
| Timestamp | 2023-05-31 12:50:45 |
| Last change | 2023-05-31 12:54:20 |

🗑 Delete   ✏ Edit   Export as Event

## Objects

Show 10 ⬍ entries     Search: _____

| Type | | Id | Tags | |
|---|---|---|---|---|
| ⬤ onion | | jamblery7zgxknhjtmj3mhfdajmyddqxbufrf6voa32h5w4otux3crqd.onion | infoleak:automatic-detection="onion"   infoleak:automatic-detection="pgp-public-key-block" | 🗑 |
| ⬤ onion | | bitmixhft4cpnciuhwffussk23tvowswbe4ttrdree74oxjmz2vyqqd.onion | infoleak:automatic-detection="onion" | 🗑 |
| 🔑 | key | 0xD3B280956F0E7CAF | | 🗑 |
| ⬤ | mail | support@jambler.io | | 🗑 |
| ⬤ | telegram | jambler | | 🗑 |
| ⬤ | name | Jambler.io | | 🗑 |

Showing 1 to 6 of 6 entries     Previous   1   Next

**TLP:CLEAR**

## AIL Framework: Extensible Capabilities

- Extending AIL to add a new **analysis module** can be done in 50 lines of Python.
- The framework **supports multi-processors/cores by default**. Any analysis module can be started multiple times to support faster processing during peak times or bulk import.
- **Multiple** concurrent **data inputs**.
- Tor Crawler (handles cookies and authentication).
- Feeders: Discord, Telegram, ...

## Ongoing Developments

- **Mail Search – Next Release**
- **Lacus crawler improvements:** proxy selection, local cache, and cookie state transfer
- **Translate search:** support for searching in other languages
- **Advanced video processing and extraction**
- **MISP export with new correlation types**
- **Automatic geolocation**

# Links

- AIL project website `https://ail-project.org`
- AIL project open source framework `https://github.com/ail-project`
- Training materials `https://github.com/ail-project/ail-training`
- Online chat `https://gitter.im/ail-project/community`



ail project

## Thank you for your attention

- AIL project[7] : https://github.com/ail-project/ail-framework
- For questions, contact: info@circl.lu

---

[7]All techniques and indicators mentioned in these slides are implemented in the AIL project, using an instance backed by a three-year dataset collected from Tor hidden services and various social networks.