

MISP Web Scraper

CTIS - October 19/20 - 2022

cudeso.be
We Secure You

<https://www.cudeso.be>
koen.vanimpe@cudeso.be

Koen Van Impe

- **Freelancer**
 - Incident response, threat intelligence, security monitoring
 - **Open source contributions**
 - MISP modules, taxonomies, automation and integration with DFIR tools, ...
 - “*MISP tip-of-the-week*”
 - **BelgoMISP**
 - Belgian MISP User Group
 - **OSINT threat feed**
 - botvrij.eu
- koen.vanimpe@cudeso.be

<https://www.cudeso.be>

<https://www.vanimpe.eu>

<https://github.com/cudeso>

@cudeso



Why?

How it started: adding threat events to botvrij.eu

Public reports Blog posts Whitepapers “OSINT”

osint:source-type="blog-post" x tlp:white x type:OSINT x



- A lot of **copy and pasting**
 - Tedious
 - Time consuming
 - Risk of errors

THE DFIR REPORT
Real Intrusions by Real Attackers, The Truth Behind the Intrusion

ANALYSTS About Unit 42 Services Threat Research Partners

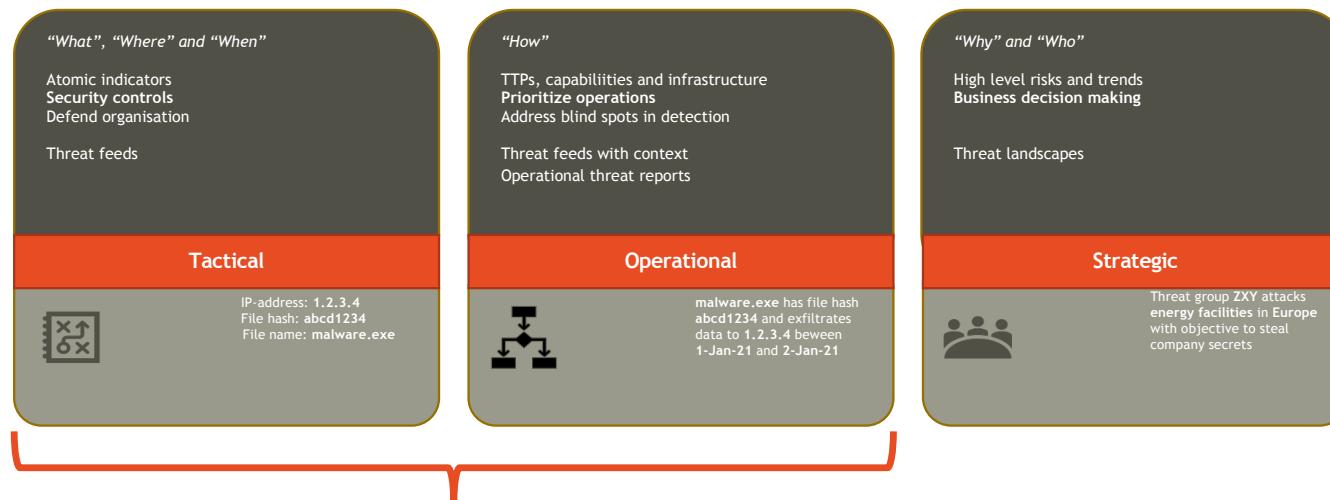
cisco TALOS Software Vulnerability Information Reputation Center Library Support Incident

SECURELIST by Kaspersky

Solutions Industries Products Services

Avoid being just an indicator list

- Focus on **tactical** and **operational** intelligence
- Include context for events and attributes

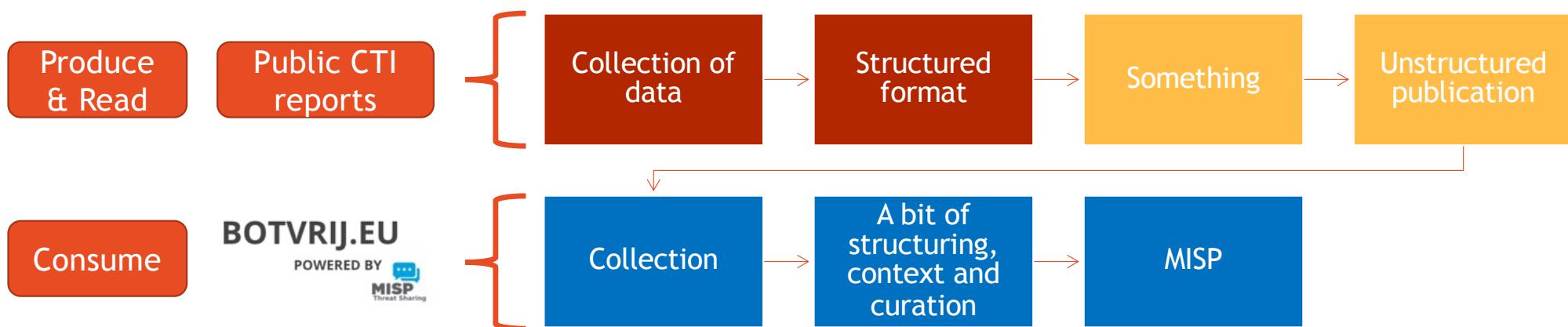
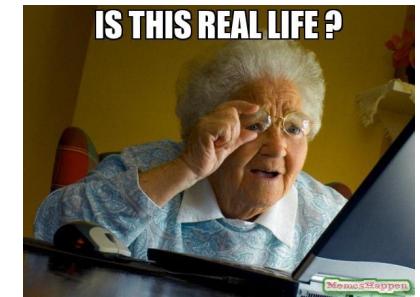


BOTVRIJ.EU

POWERED BY
 MISP
Threat Sharing

Sources used for the OSINT feed

- A lot of public reports *mostly contain unstructured information*
 - Or structured but not consistent in format
 - CSV in format x, CSV in format y, TXT files
- Reports sometimes **focus more on readers**, than on consumers
- Workflow (presumably not just for me ...)



Why are public sources so difficult to consume?

- It's not about *solving the sharing problem*
 - But about making public reports more easily consumable
- There are commercial interests
 - But providing structured information publicly can be a valuable “teaser” for your professional CTI services
- Not all is bad

WHY REINVENT THE WHEEL WHEN YOU DON'T HAVE TO?



The image displays three GitHub repository pages side-by-side, illustrating the practice of "teasing" through public commits:

- eset / malware-ioc**: Shows contributions from users like porolli, amavaldo, animalfarm, attor, and backdoordiplomacy, all adding IoCs for various threat actors.
- pan-unit42 / iocs**: Shows contributions from chkroot, AcidBox, ArtraDownloader, and BazarCall, all adding IoCs for Qakbot and other malware families.
- NVISOsecurity / nviso-cti**: Shows contributions from sylktools, adding IoCs for the Iranian Threat Actor and bumblebeeRoundTwo, which are described as being used for September's Generic Threat Briefing.

What would make this work easier?

- Formats?
 - STIX
 - Low adoption
 - Certainly what concerns public reports
 - CSV/TXT
 - Not always consistent
 - Field sequence, naming conventions
 - PDF/XLS/...
 - No



What would make this work easier?

- Formats?
 - STIX
 - Low adoption
 - Certainly what concerns public reports
 - CSV/TXT
 - Not always consistent
 - Field sequence, naming conventions
 - PDF/XLS/...
 - No
 - Pending everyone starts using the MISP core format ...



What would make this work easier?

- Formats?

- STIX

- Low adoption
 - Certainly what concerns put

- CSV/TXT

- Not always consistent
 - Field sequence, naming con\

- PDF/XLS/...

- No

- Pending everyone starts using the **MISP core format** ...

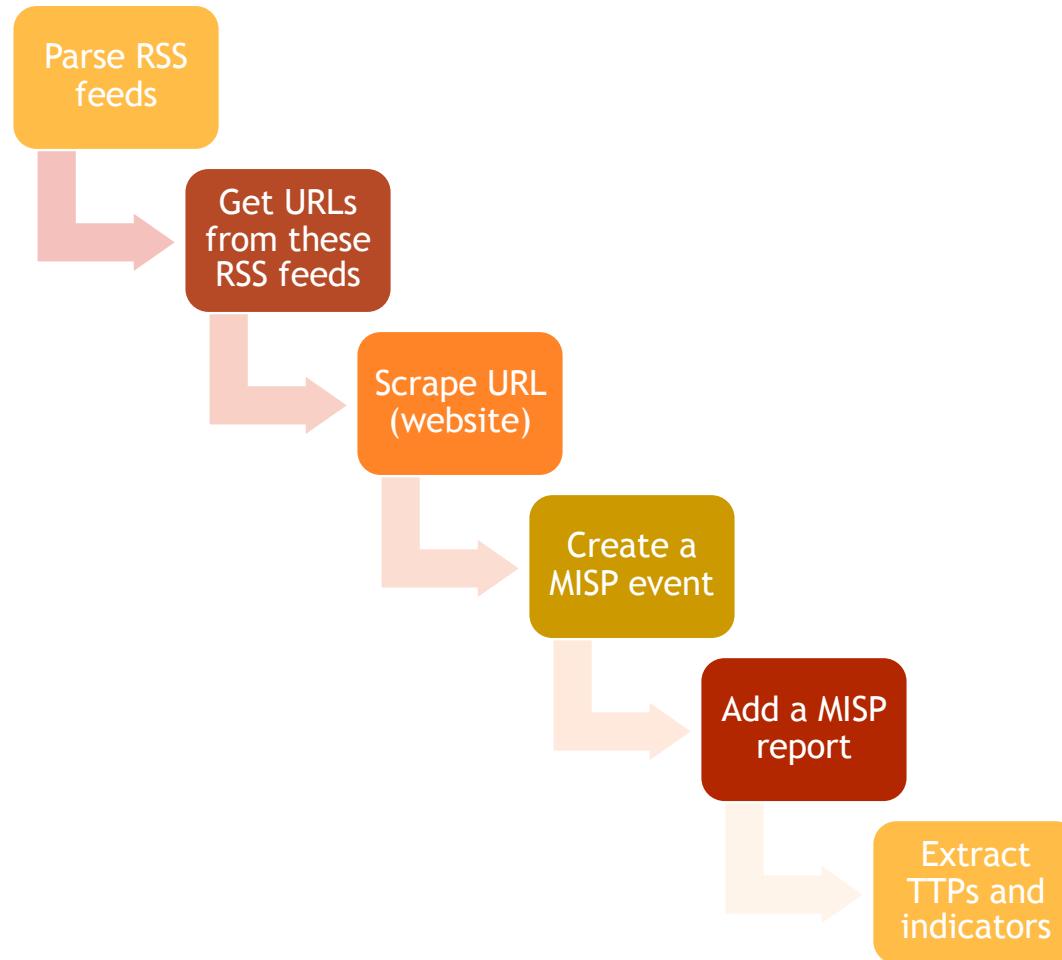
The screenshot shows a GitHub repository page for 'eset / malware-ioc'. The repository is public. The 'Issues' tab is selected, showing one open issue. The issue is titled 'Create a MISP shareable files #20'. A comment from 'ZLT-ops' is visible, dated 3 days ago. The comment text is:

ZLT-ops commented 3 days ago

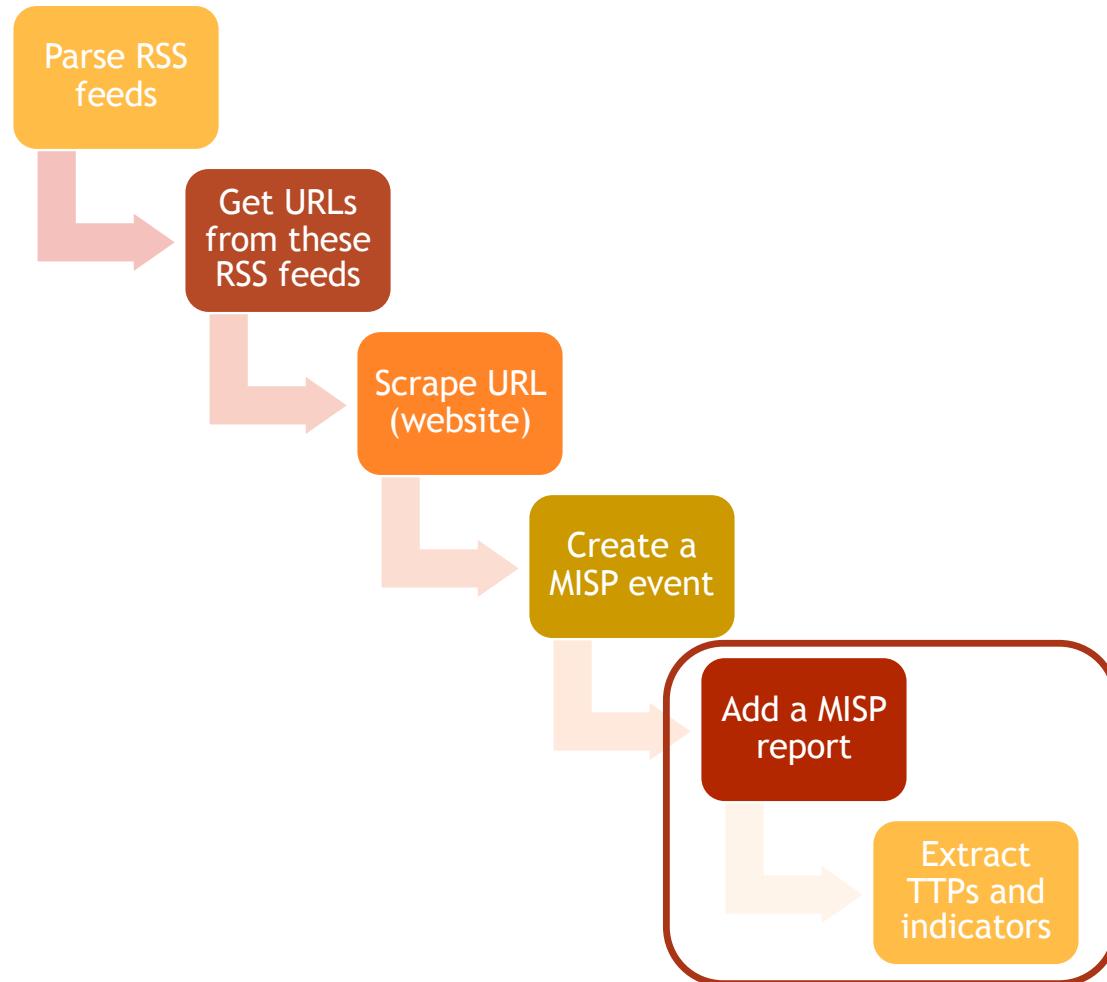
Hi Team,
Could you share a MISP ready file for your repo so users could just import the IOCs to their MISP?
Thanks!

MISP web scraper

MISP web scraper



MISP web scraper



Build on MISP reports

"CTI report"

References to
attributes and
objects

References to
taxonomies and
galaxy matrices

In raw text or
markdown

Is part of an
event and
distributed with
events

Build on MISP reports

"CTI report"

References to
attributes and
objects

References to
taxonomies and
galaxy matrices

In raw text or
markdown

Is part of an
event and
distributed with
events

1

Indicators

TTPs

Build on MISP reports

"CTI report"

References to attributes and objects

References to taxonomies and galaxy matrices

In raw text or markdown

Is part of an event and distributed with events

1

Indicators

TTPs

2

Convert "website" content to markdown

main ▾ [misp-modules / misp_modules / modules / expansion / html_to_markdown.py](#)

Event Reports

Import from URL

Build on MISP reports

"CTI report"

References to attributes and objects

References to taxonomies and galaxy matrices

In raw text or markdown

Is part of an event and distributed with events

1

Indicators

TTPs

2

Convert "website" content to markdown

main ▾ misp-modules / misp_modules / modules / expansion / html_to_markdown.py

Event Reports

Import from URL

3

Save Menu ▾ Help

Download

Markdown parsing rules

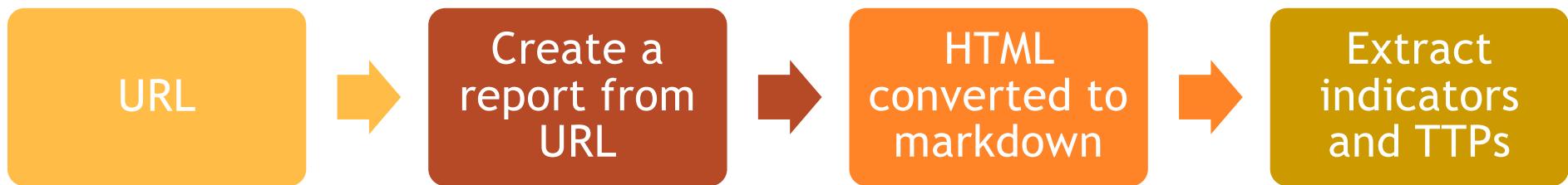
Markdown rendering rules

Extract entities

Manual extraction

Automatic extraction

Easy wins with the MISP report feature



MISP web scraper components

Flask

- Web form
- Manual: add a URL to Redis

MISP-Scraper

Go to MISP

Add either a page title and URL or raw HTML to submit to MISP Scraper

Feed title

Manual

Feed URL

Manual

Page title

Leave blank to use the HTML title

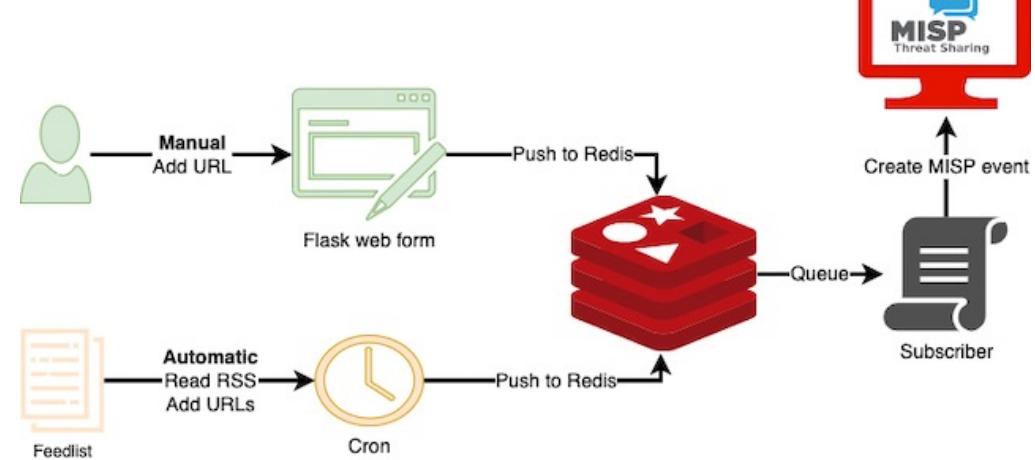
Page URL (link)

(not scraped when raw HTML is submitted)

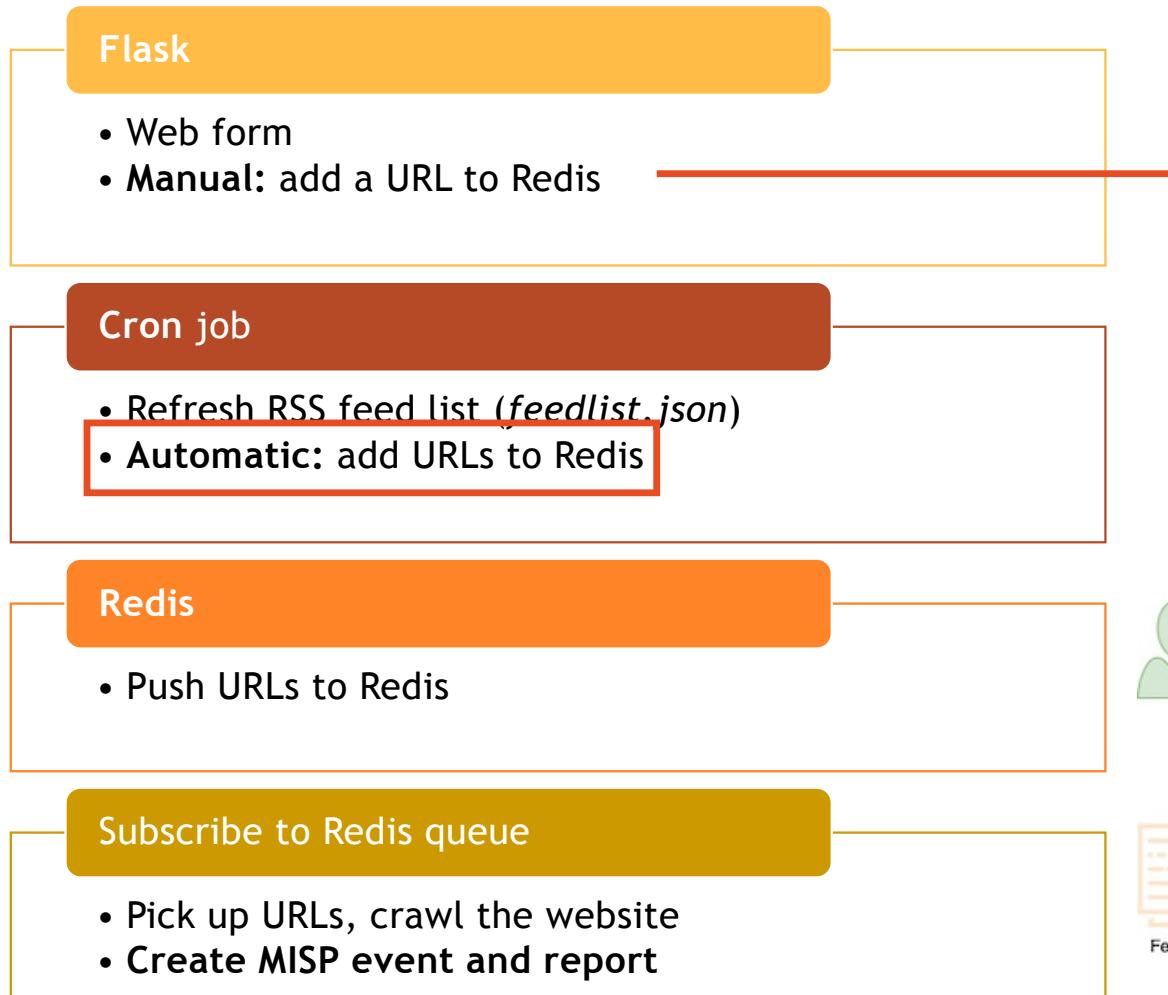
URL to scrape

Raw HTML

(ignored when empty, then Page URL is scraped)



MISP web scraper components



MISP-Scraper [Go to MISP](#)

Add either a page title and URL or raw HTML to submit to MISP Scraper

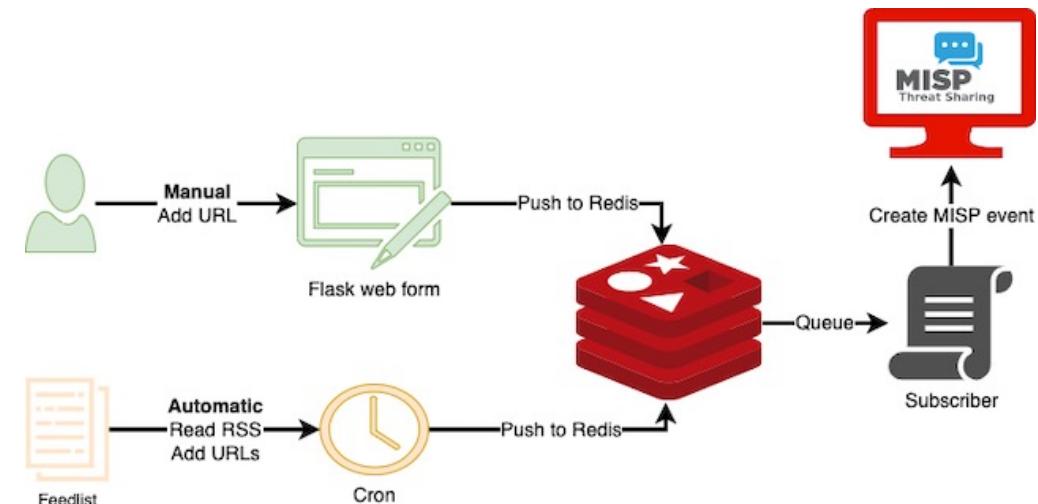
Feed title

Feed URL

Page title

Page URL (link)
(not scraped when raw HTML is submitted)

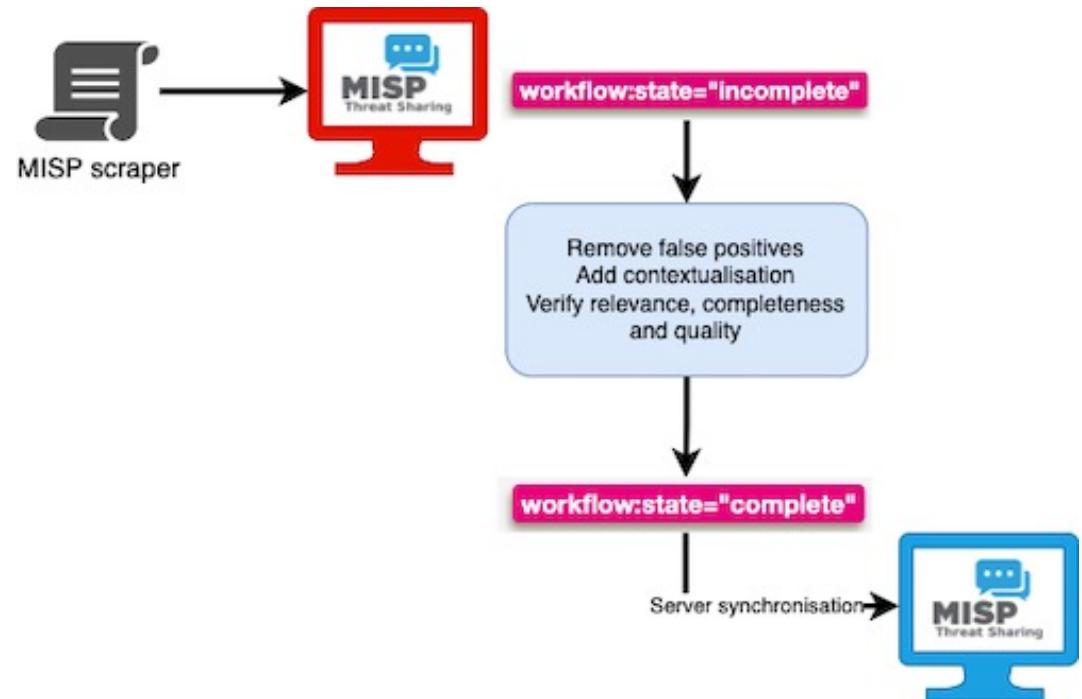
Raw HTML
(ignored when empty, then Page URL is scraped)



Web scraper workflow

- **Scrape** in a “dirty” environment
 - Automated, “cron job”
 - Create MISP event

 **workflow:state="incomplete"** x

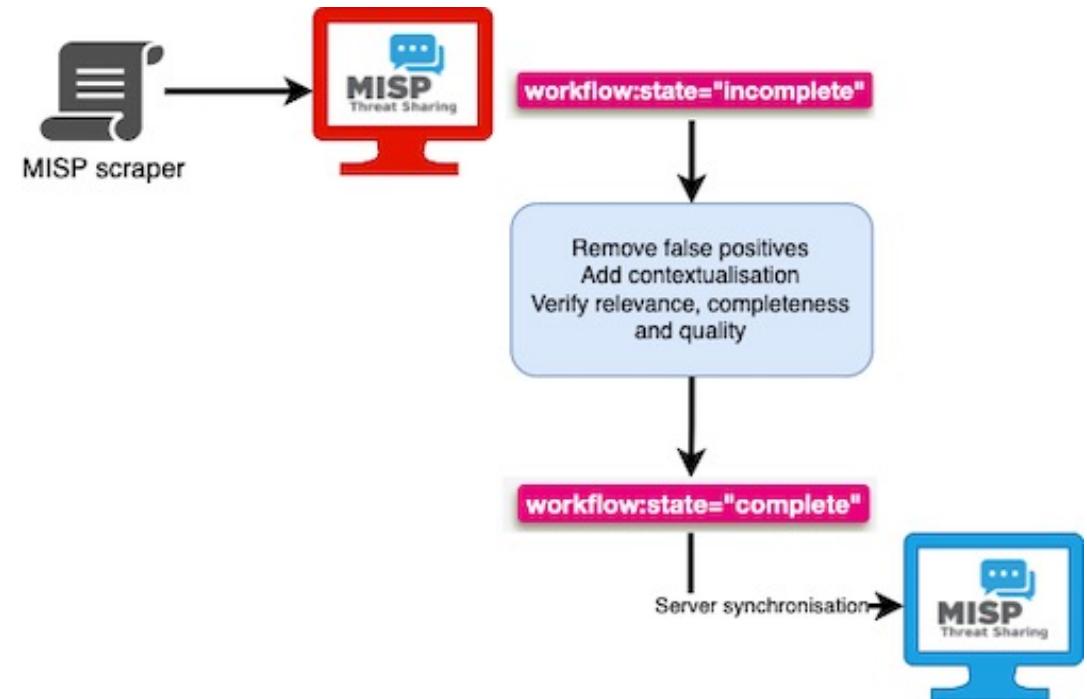


Web scraper workflow

- **Scrape** in a “dirty” environment

- Automated, “cron job”
- Create MISP event

 **workflow:state="incomplete"** x



- **Curate** event

- Scrape successful?
- Relevant and complete?
- Remove false positives

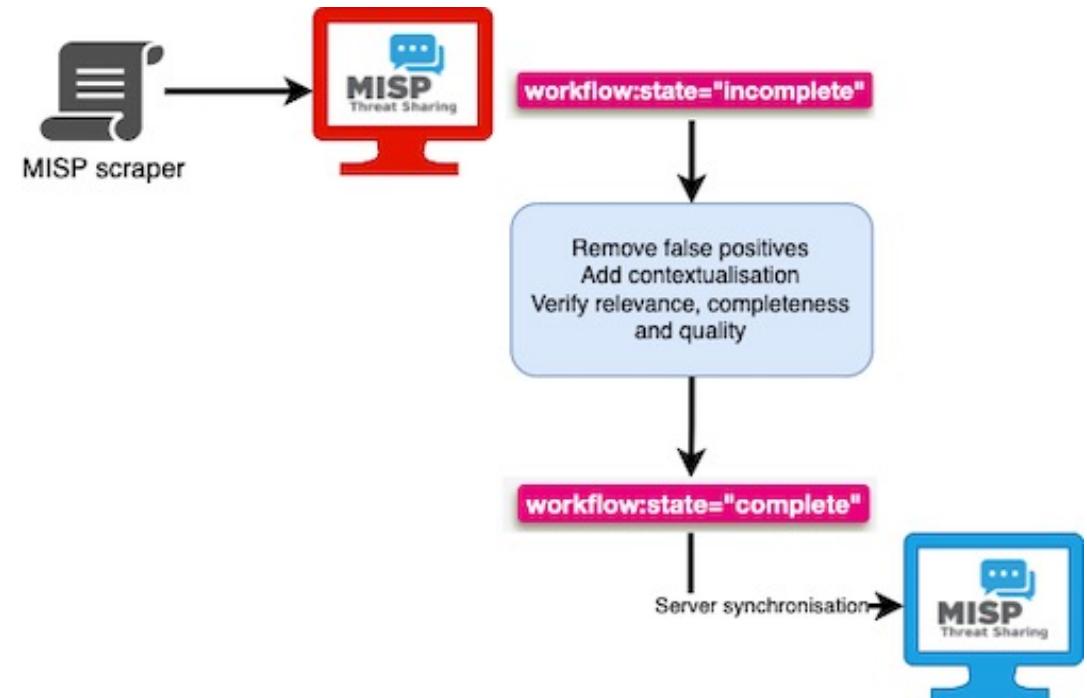
 **workflow:state="complete"** x

Web scraper workflow

- **Scrape** in a “dirty” environment

- Automated, “cron job”
- Create MISP event

 **workflow:state="incomplete"** x



- **Curate** event

- Scrape successful?
- Relevant and complete?
- Remove false positives

 **workflow:state="complete"** x

- **Publish and synchronise**

- With your “clean” MISP

Curate scraped events (1/2)

- Analyst checks
 - Scrape successful?
 - HTTP error codes added as tags
 - Not “blocked” by a CDN

misp-scaper:HTTP=403
misp-scaper:HTTP=404

html .st0{display:none}.st1{fill:#006db6} # Access denied
You cannot access **hostname**. Refresh the page or contact the site owner to request access.
• Ray ID: **[REDACTED]**
• Timestamp: 2022-10-12 08:30:40 UTC
• Your IP address: **ip-dst [REDACTED]**
• Requested URL: @[suggestion](**hostname**) /2022/10/threat-roundup-0930-1007.html
• Error reference number: 1020
• **mitre-attack-pattern** → Server - T1583.004 ID: **regkey FL__78F89**
• User-Agent: **url python-requests/2.27.1**
Cloudflare Ray ID: **[REDACTED]** Your IP: **ip-dst [REDACTED]**

<input type="checkbox"/>	Date	Org	Category	Type	Value
<input type="checkbox"/>	2022-10-12		Other	comment	Preventing Cryptocurrency Cyber Extortion
<input type="checkbox"/>	2022-10-12		External analysis	link	http://feeds.trendmicro.com/TrendMicroSimplySecurity
<input type="checkbox"/>	2022-10-12		External analysis	link	https://www.trendmicro.com/en_us/ciso/22/i/prevent-cyber-extortion.html 

Page 1 of 1, showing 1 records out of 3 total, starting on record 1, ending on 3

- Wrong attribute category or types

Payload delivery hostname 2022-09-23-iocs-for-icedid-and-cobalt-strike.txt.zip

Curate scraped events (2/2)

- **False positives**
 - MISP warninglists
- **Mass edit attributes**
 - Remove to_ids
 - Attributes primarily fit for context, not for detection or filtering
 - Remove correlation
- **Mass delete attributes**
 - Attributes irrelevant for the threat event

Mass Edit Attributes

Distribution

Do not alter current settings

For Intrusion Detection System

No

Create proposals

Correlations

Disable correlations

Payload delivery	filename	Trojan.BAT.ADFIND.YECGUT
Payload delivery	filename	Backdoor.Win32.SYSTEMBC.YXCF LZ
Payload delivery	filename	Backdoor.Win32.COBEACON.YXCH3
Payload delivery	filename	Ransom.Win32.PLAYCRYPT.YECGUT
Payload delivery	filename	Trojan.Win64.PRIVICMD.YXCHW
Payload delivery	filename	Ransom.Win32.PLAYDE.YACHWT

Make curation a bit easier (1)

- **Doubles**

- Web scraper: skip MISP events already created
 - (based on event title + scraper tag)

- **Retention**

- Old events are automatically deleted
- Avoid analyst overload

retention:2d

retention:7d

- **Classification of scraped events**



- Tool and source
 - Unique tag per source
 - Allows f.e. to filter for all posts scraped from the source TrendMicro
- Default TLP and workflow state

Make curation a bit easier (2)

- Custom warninglist
 - Via a MISP warninglist
 - Avoid adding known non-relevant attributes
 - URLs to partner websites
 - Social media
 - “Commercial” strings
- Links and comment attributes to describe source
 - Feed, title
 - No correlation
 - Link to web page
 - Correlation

SCRAPER

ID	72
Name	Scraper
Description	Warninglist to be used for the web scraper
Version	9
Category	False positive
Type	string
Accepted attribute types	
Enabled	Yes Disable

Values

<http://club.orbisius.com/products/wordpress-plugins/whitelist-ip-for-limit->
<https://t.co/WpmpTPS2T>
<https://threatpost.com/malware-traffic-analysis.net/payloads.ThisZone.Identifier>
[src="/static/blogv2/js/vendors.js](src=)
<http://google.com/ads/remarketingsetup>
<http://threatpost.com>

External analysis	link	https://securelist.com/feed/ ⚠	🌐 + 👤 + 🌐 + 👤 +	Feed URL <input type="checkbox"/>	🔍
External analysis	link	https://securelist.com/malicious-whatsapp-mod-distributed-through-legitimate-apps/107690/ ⚠	🌐 + 👤 + 🌐 + 👤 +	Blog URL <input checked="" type="checkbox"/>	🔍
Other	comment	Malicious WhatsApp mod distributed through legitimate apps	🌐 + 👤 + 🌐 + 👤 +	Blog title <input type="checkbox"/>	🔍

Make curation a bit easier (2)

- Custom warninglist
 - Via a MISP warninglist
 - Avoid adding known non-relevant attributes
 - URLs to partner websites
 - Social media
 - “Commercial” strings

Curation: events with less than 4 attributes can be ignored

- Links and comment attributes to describe source report
 - Feed, title
 - No correlation
 - Link to web page
 - Correlation

SCRAPER

ID	72
Name	Scraper
Description	Warninglist to be used for the web scraper
Version	9
Category	False positive
Type	string
Accepted attribute types	
Enabled	Yes Disable

Values

<http://club.orbisius.com/products/wordpress-plugins/whitelist-ip-for-limit-l>
<https://t.co/VWpmpTPS2T>
<https://threatpost.com/malware-traffic-analysis.net/payloads.This.Zone.Identifier>
[src="/static/blogv2/js/vendors.js](src=)
<http://google.com/ads/remarketingsetup>
<http://threatpost.com>

External analysis	link	https://securelist.com/feed/ ⚠	🔗 👤 🔗 👤	Feed URL	<input type="checkbox"/>	🔍
External analysis	link	https://securelist.com/malicious-whatsapp-mod-distributed-through-legitimate-apps/107690/ ⚠	🔗 👤 🔗 👤	Blog URL	<input checked="" type="checkbox"/>	🔍
Other	comment	Malicious WhatsApp mod distributed through legitimate apps	🔗 👤 🔗 👤	Blog title	<input type="checkbox"/>	🔍

Scraped event

Preventing Black Basta Ransomware in 2022

Event ID	294
UUID	391357f3-d319-41f8-8525-ce2b083097b6 +≡
Creator org	CUDESO
Owner org	CUDESO
Creator user	[REDACTED]
Protected Event (experimental) ⓘ	Tags from scraper to ease curation
Tags	 misp:tool="misp-scra..."/> osint:source-type="blog-post"/> misp:event-type="collection"/> misp-galaxy:botnet="Qbot"/> misp-galaxy:malpedia="Black Basta"/> workflow:state="complete"/> tip:white + +
Date	2022-09-12

Scraped event

Preventing Black Basta Ransomware in 2022

Event ID 294

UUID 391357f3-d319-41f8-8525-ce2b083097b6

Creator org CUDESO

Owner org CUDESO

Creator user [REDACTED]

Protected Event (experimental) Tags from scraper to ease curation

Tags

Date 2022-09-12

Attribute	Type	Value	Actions
2022-09-12	Network activity	url https://aazsbssgya565vlu2c6bzy6yfiebkcbtvvcyvtolt33s77xyp7nypyd.onion/	Inherit
2022-09-12	Payload delivery	sha256 9fce9ee85516533bae34fc1184a7cf31fa9f2c7889b13774f83d1df561708833	Inherit
2022-09-12	Payload delivery	sha256 203d2807df6ef531efbec7bfd109986de3e23df64c01ea4e337cbe5ba675248b	Inherit
2022-09-12	Payload delivery	sha256 0d6c3de5aebbbe85939d7588150edf7b7bcd712fc6a83d79e65b6f79bf2ef	Inherit
2022-09-12	Payload delivery	sha1 2bee3f716b80273db9639376a296cf19cdba0f1a	Inherit
2022-09-12	Network activity	url http://146.70.79.52/	Inherit

MISP report

Event Reports

+ Add Event Report Import from URL Generate report from Event All Default Deleted

ID	Name	Last update	Distribution	Actions
16	Report from - https://www.deepinstinct.com/blog/black-basta-ransomware-threat-emergence (1662992420)	2022-09-12 16:20:20	Inherit event	

Extracted attributes

Scraped event

Ragnar Locker Ransomware Targeting the Energy Sector

Event ID	292
UUID	8dbeaac-a671-4a02-8dab-5eec2a1c935b  
Creator org	CUDESO
Owner org	CUDESO
Creator user	koen.vanimpe@cudeso.be
Protected Event (experimental) 	 Event is in unprotected mode.
Tags	 misp:tool="misp-scaper"   osint:source-type="blog-post"   misp:event-type="collection"   misp-galaxy:malpedia="RagnarLocker (Windows)"   workflow:state="complete"   tip:white   
Date	2022-09-12

Scraped event

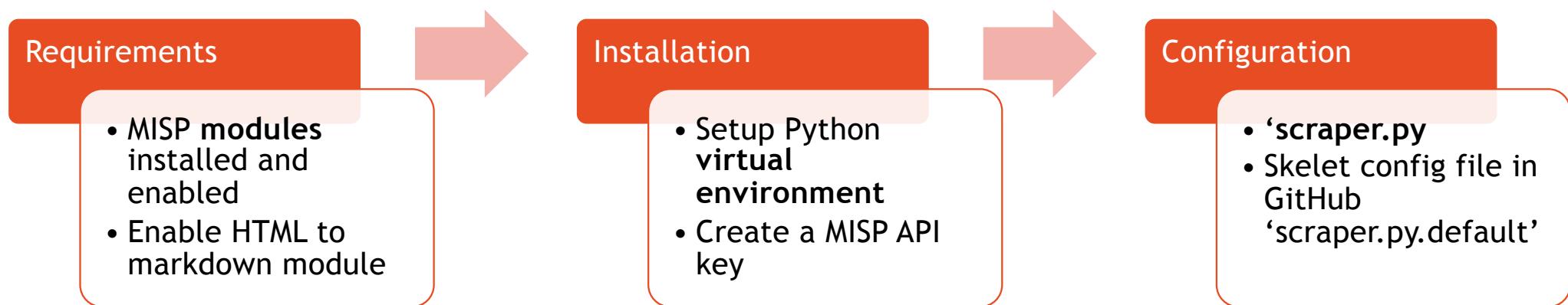
Ragnar Locker Ransomware Targeting the Energy Sector

Event ID	292																																														
UUID	8dbeaac-a671-4a02-8dab-5eec2a1c935b																																														
Creator org	CUDESO																																														
Owner org	CUDESO																																														
Creator user	koen.vanimpe@cudeso.be																																														
Protected Event (experimental)	Event is in unprotected mode.																																														
Tags																																															
Date	2022-09-12																																														
	<table><tbody><tr><td><input type="checkbox"/> 2022-09-12 Payload delivery filename ntuser.dat.log</td><td> </td><td> </td><td><input checked="" type="checkbox"/></td><td><input type="checkbox"/></td></tr><tr><td><input type="checkbox"/> 2022-09-12 Payload delivery filename bootfront.bin</td><td> </td><td> </td><td><input checked="" type="checkbox"/></td><td><input type="checkbox"/></td></tr><tr><td><input type="checkbox"/> 2022-09-12 Payload delivery sha256 9b62cdb57f4c34924333dfa3baef9d93eefab68109580b682b0740fe73b63983</td><td> </td><td> </td><td> <input checked="" type="checkbox"/></td><td><input checked="" type="checkbox"/></td></tr><tr><td><input type="checkbox"/> 2022-09-12 Payload delivery sha256 c2bd70495630ed8279de0713a010e5e55f3da29323b59ef71401b12942ba5216</td><td> </td><td> </td><td> <input checked="" type="checkbox"/></td><td><input checked="" type="checkbox"/></td></tr><tr><td><input type="checkbox"/> 2022-09-12 Payload delivery sha256 5469182495d92a5718e0e1cdf371e92b79724e427050154f318e693d341c89</td><td> </td><td> </td><td> <input checked="" type="checkbox"/></td><td><input checked="" type="checkbox"/></td></tr><tr><td><input type="checkbox"/> 2022-09-12 Payload delivery sha256 ec35c76ad2c8192f09c02eca1f263b406163470ca8438d054db7adcf5bf0597</td><td> </td><td> </td><td> <input checked="" type="checkbox"/></td><td><input checked="" type="checkbox"/></td></tr><tr><td><input type="checkbox"/> 2022-09-12 Payload delivery sha256 dd5d4cf9422b6e4514d49a3ec542cffb682be8a24079010cda689afbb44ac0f4</td><td> </td><td> </td><td> <input checked="" type="checkbox"/></td><td><input checked="" type="checkbox"/></td></tr><tr><td><input type="checkbox"/> 2022-09-12 Payload delivery sha256 cf5ec678a2f836f859eb983eb633d529c25771b3b7505e74aa695b7ca00f9fa8</td><td> </td><td> </td><td> <input checked="" type="checkbox"/></td><td><input checked="" type="checkbox"/></td></tr><tr><td><input type="checkbox"/> 2022-09-12 Payload delivery sha256 ce33096639fb5c51684e9e3a7c7c161884ecad29e8d6ad602fd8be42076b8d4</td><td> </td><td> </td><td> <input checked="" type="checkbox"/></td><td><input checked="" type="checkbox"/></td></tr></tbody></table>	<input type="checkbox"/> 2022-09-12 Payload delivery filename ntuser.dat.log			<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> 2022-09-12 Payload delivery filename bootfront.bin			<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> 2022-09-12 Payload delivery sha256 9b62cdb57f4c34924333dfa3baef9d93eefab68109580b682b0740fe73b63983			<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/> 2022-09-12 Payload delivery sha256 c2bd70495630ed8279de0713a010e5e55f3da29323b59ef71401b12942ba5216			<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/> 2022-09-12 Payload delivery sha256 5469182495d92a5718e0e1cdf371e92b79724e427050154f318e693d341c89			<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/> 2022-09-12 Payload delivery sha256 ec35c76ad2c8192f09c02eca1f263b406163470ca8438d054db7adcf5bf0597			<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/> 2022-09-12 Payload delivery sha256 dd5d4cf9422b6e4514d49a3ec542cffb682be8a24079010cda689afbb44ac0f4			<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/> 2022-09-12 Payload delivery sha256 cf5ec678a2f836f859eb983eb633d529c25771b3b7505e74aa695b7ca00f9fa8			<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/> 2022-09-12 Payload delivery sha256 ce33096639fb5c51684e9e3a7c7c161884ecad29e8d6ad602fd8be42076b8d4			<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	
<input type="checkbox"/> 2022-09-12 Payload delivery filename ntuser.dat.log			<input checked="" type="checkbox"/>	<input type="checkbox"/>																																											
<input type="checkbox"/> 2022-09-12 Payload delivery filename bootfront.bin			<input checked="" type="checkbox"/>	<input type="checkbox"/>																																											
<input type="checkbox"/> 2022-09-12 Payload delivery sha256 9b62cdb57f4c34924333dfa3baef9d93eefab68109580b682b0740fe73b63983			<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>																																											
<input type="checkbox"/> 2022-09-12 Payload delivery sha256 c2bd70495630ed8279de0713a010e5e55f3da29323b59ef71401b12942ba5216			<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>																																											
<input type="checkbox"/> 2022-09-12 Payload delivery sha256 5469182495d92a5718e0e1cdf371e92b79724e427050154f318e693d341c89			<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>																																											
<input type="checkbox"/> 2022-09-12 Payload delivery sha256 ec35c76ad2c8192f09c02eca1f263b406163470ca8438d054db7adcf5bf0597			<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>																																											
<input type="checkbox"/> 2022-09-12 Payload delivery sha256 dd5d4cf9422b6e4514d49a3ec542cffb682be8a24079010cda689afbb44ac0f4			<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>																																											
<input type="checkbox"/> 2022-09-12 Payload delivery sha256 cf5ec678a2f836f859eb983eb633d529c25771b3b7505e74aa695b7ca00f9fa8			<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>																																											
<input type="checkbox"/> 2022-09-12 Payload delivery sha256 ce33096639fb5c51684e9e3a7c7c161884ecad29e8d6ad602fd8be42076b8d4			<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>																																											
		<table><thead><tr><th colspan="2">Event Reports</th><th>Last update</th><th>Distribution</th><th>Actions</th></tr><tr><th>ID</th><th>Name</th><th></th><th></th><th></th></tr></thead><tbody><tr><td>14</td><td>Report from - https://www.cybereason.com/blog/threat-analysis-report-ragnar-locker-ransomware-targeting-the-energy-sector (1662984886)</td><td>2022-09-12 14:14:47</td><td>Inherit event</td><td></td></tr></tbody></table>	Event Reports		Last update	Distribution	Actions	ID	Name				14	Report from - https://www.cybereason.com/blog/threat-analysis-report-ragnar-locker-ransomware-targeting-the-energy-sector (1662984886)	2022-09-12 14:14:47	Inherit event																															
Event Reports		Last update	Distribution	Actions																																											
ID	Name																																														
14	Report from - https://www.cybereason.com/blog/threat-analysis-report-ragnar-locker-ransomware-targeting-the-energy-sector (1662984886)	2022-09-12 14:14:47	Inherit event																																												

What works well?

Always	Often	Sometimes
<ul style="list-style-type: none">• Hashes• IP address• Easy to identify with regexp	<ul style="list-style-type: none">• URLs• Domain and hostname• E-mail src• TTPs<ul style="list-style-type: none">• If not part of the commercial text on web pages	<ul style="list-style-type: none">• Filenames• Registry keys
<ul style="list-style-type: none">• Good for indicator collection• Some automatic TTPs• Converting to objects remains a manual process• Detection rules (yara) require manual edition• RSS reader for threat reports		

Installation and configuration



```
git clone https://github.com/cudeso/misp-scraper  
cd misp-scraper  
virtualenv scraper  
source scraper/bin/activate  
pip install -r requirements.txt  
cp scraper.py.default scraper.py
```

Configuration file

```
# General configuration
misp_scraper_log = "/var/log/misp/scraper.log"                                # Log file for the scraper
app_name = "MISP-Scraper"                                                       # Used in the Flask server
manual_feed = "Manual"
manual_feedsource = "Manual"
flask_secret_key = "MyZsuePerSecret1984key"                                     # Flask secret key. Change this to
flask_address = "127.0.0.1"                                                       # IP address for the Flask server
flask_port = 5200
flask_certificate_file = "/etc/ssl/private/misp.local.crt"                      # Point to a certificate file (fe.
flask_certificate_keyfile = "/etc/ssl/private/misp.local.key"                   # Point to a certificate key file
logging_level = "error"
feedlist = "/home/ubuntu/misp-scraper/feedlist.json"                            # Location of the RSS feedlist
rawhtml_distribution = 4
rawhtml_sharing_group_id = 4

# MISP configuration
misp_key = ""                                                                    # MISP API key
misp_url = ""                                                                    # MISP url
misp_verifycert = False
misp_distribution = 2
misp_threat_level_id = 2
misp_analysis_level = 2
misp_scraper_event = "Scraper"
misp_retentiontime = "99d"
misp_scraper_tags = ["misp:tool=\"misp-scraper\"", "osint:source-type=\"blog-post\"", "misp:event-"
misp_warninglist = 0
misp_hard_delete_on_cleanup = False
```

Feedlist configuration file



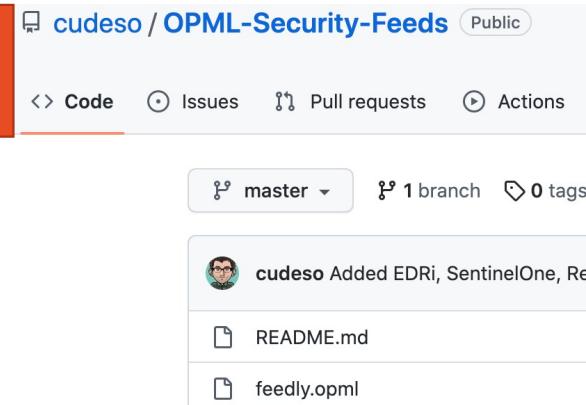
Feedlist configuration file

JSON

List of feeds

- title
 - Used in the tagging of events
- url
 - URL to download the RSS feed

- Need inspiration for RSS feeds?
 - <https://github.com/cudeso/OPML-Security-Feeds>



Room for improvement (1)

- Web crawler is Python
 - User-agent blocked
 - Improve with **ail-splash-manager**
 - Instead of handling the crawling in web-scraping, have AIL do the crawling
 - Microsoft Playwright ^^AIL

The screenshot shows a web browser displaying a 403 Access Denied error page. The page contains the following elements:

- MISP Metadata:** misp:tool="misp-scaper", osint:source-type="blog-post", retention:2d, misp:event-type="collection", tlp:white, workflow:state="Incomplete", scraper:, misp-scaper:HTTP=403.
- Timeline:** Ray ID: [REDACTED], Timestamp: 2022-10-12 08:30:40 UTC, Your IP address: ip-dst [REDACTED], Requested URL: @[suggestion](hostname [REDACTED]) / 2022/10/threat-roundup-0930-1007.html, Error reference number: 1020, mitre-attack-pattern => Server - T1583.004 ID: regkey FL_78F89, User-Agent: url python-requests/2.27.1.
- Cloudflare Information:** Cloudflare Ray ID: [REDACTED], Your IP: ip-dst [REDACTED].
- Navigation:** main ▾ DocIntel / DocIntel.Services.Scrapers /
- Footer:** Antoine Cailliau Release for CTI Summit

Room for improvement (1)

- Web crawler is Python
 - User-agent blocked
 - Improve with `ail-splash-manager`
 - Instead of handling the crawling in web-scraping, have AIL do the crawling
 - Microsoft Playwright ^^AIL
 - First steps done by allowing to submit raw HTML
- Include screenshots
 - Via `ail-splash-manager`?
- Deal with RSS sources containing ridiculous number of URLs
 - Feeds containing all articles going back to 20xx
 - Discover publish date and only consider recent reports

The screenshot shows a MISP event page with several tags: misp:tool="misp-scaper", osint:source-type="blog-post", retention:2d, misp:event-type="collection", tlp:white, workflow:state="Incomplete", scraper:, and misp-scaper:HTTP=403. Below the tags, there is a message: "html .st0[display:none],st1{fill:#006db6} # Access denied You cannot access hostname. Refresh the page or contact the site owner to request access." A list of details follows: Ray ID, Timestamp (2022-10-12 08:30:40 UTC), Your IP address (ip-dst), Requested URL (@[suggestion](hostnam 2022/10/threat-roundup-0930-1007.html), Error reference number (1020), mitre-attack-pattern (Server - T1583.004 ID: regkey FL_78F89), and User-Agent (url python-requests/2.27.1). At the bottom, it says Cloudflare Ray ID: [redacted] Your IP: ip-dst [redacted]. On the right, there's a navigation bar with 'main' and 'DocIntel / DocIntel.Services.Scaper /' and a note 'Antoine Cailliau Release for CTI Summit'.

Room for improvement (2)

- Use Natural Language Processing (**NLP**) to analyse reports
 - FIRSTCON 2022

All the Unstructured Data! Using NLP to Process Threat Reports

TLP:AMBER  ↗

Patrick Grau (Bosch, DE)

- Bulk change category and type of attributes
 - Similar as to bulk edit events
 - Would in general be a great feature to edit MISP events
- Use Feedly AI
 - Remove non-relevant clutter text from web pages

Room for improvement (3)

- Automatically **remove attributes** for specific scraped sources
 - F.e. always remove attribute xyz from source abc
 - Company labels etc.
 - Remove specific **CSS classes**

Room for improvement (3)

- Automatically **remove attributes** for specific scraped sources
 - F.e. always remove attribute xyz from source abc
 - Company labels etc.
 - Remove specific **CSS classes**
 - Need galaxies/clusters?

Misinformation Pattern

-  Instagram 
-  Facebook 
-  Twitter 

Target Information

-  Australia 
-  Canada 
-  France 
-  Hong Kong 
-  Indonesia 
-  Ireland 
-  Malaysia 
-  New Zealand 
-  Philippines 
-  Singapore 
-  South Africa 
-  United Kingdom 
-  United States 

Sector

-  Development 
-  Education 
-  Electric 
-  Gas 
-  Health 
-  Industrial 
-  Insurance 
-  Intelligence 
-  Legal 
-  Manufacturing 
-  Oil 
-  Technology 
-  Water 

Room for improvement (3)

- Automatically remove attributes for specific scraped sources
 - F.e. always remove attribute xyz from source abc
 - Company labels etc.
- Remove specific CSS classes
 - Need galaxies/clusters?
 - Mass delete clusters

The screenshot shows a dark-themed web interface with a sidebar containing filtering options and a main content area displaying a dropdown menu's HTML code.

Misinformation Pattern

- Instagram
- Facebook
- Twitter

Target Information

- Australia
- Canada
- France
- Hong Kong
- Indonesia
- Ireland
- Malaysia
- New Zealand
- Philippines
- Singapore
- South Africa
- United Kingdom
- United States

Sector

- Development
- Education
- Electric
- Gas
- Health
- Industrial
- Insurance
- Intelligence
- Legal
- Manufacturing
- Oil
- Technology
- Water

HTML Code (dropdown menu):

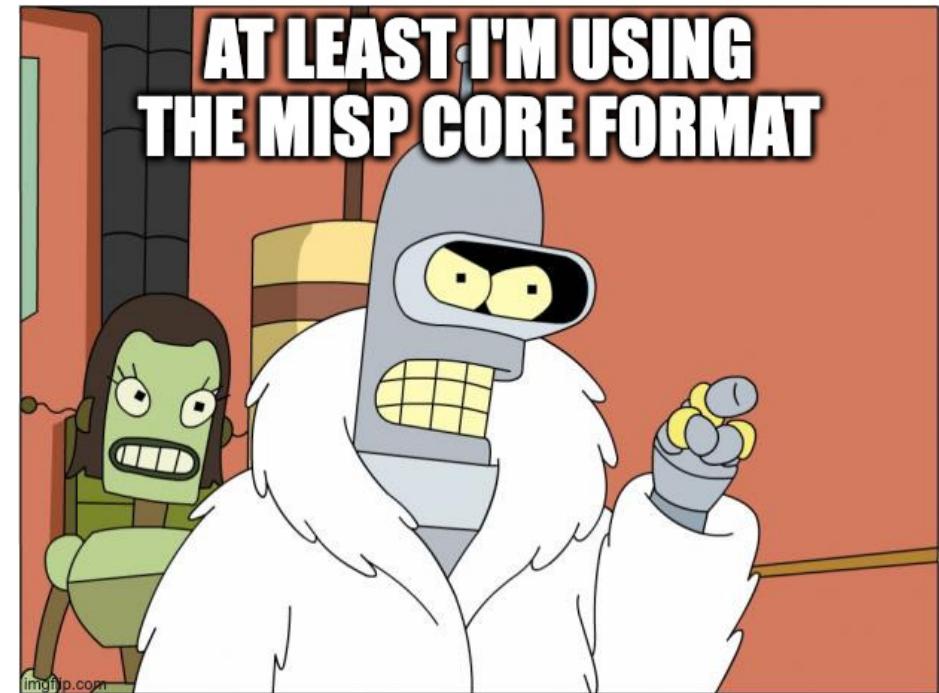
```
<div class="dropdown stretched-dropdown">
  <button class="menu-button button-default" type="button" data-toggle="dropdown">
    <span class="menu-button__icon icon-region"></span>
    <span class="menu-button__text">Region</span>
  </button>
  <div class="dropdown-menu align-left">
    <ul class="menu-column col-xs-12 col-sm-4 col-md-3">
      <li class="dropdown-header">
        The Americas
      </li>
      <li>
        <a href="/en_us.html">
          United States
        </a>
      </li>
      <li>
        <a href="/pt_br.html">
          Brasil
        </a>
      </li>
    </ul>
  </div>
</div>
```

Room for improvement (4)

- Scraping websites to gather threat intel reports is far from ideal.

Room for improvement (4)

- Scraping websites to gather threat intel reports is far from ideal.
 - MISP web scraper is no longer needed and everyone uses the MISP core format.



Resources

Get it via GitHub

- <https://github.com/cudeyo/misp-scraper>

Documentation on GitHub and on the MISP project website

- <https://www.misp-project.org/2022/08/08/MISP-scraper.html/>

Questions?