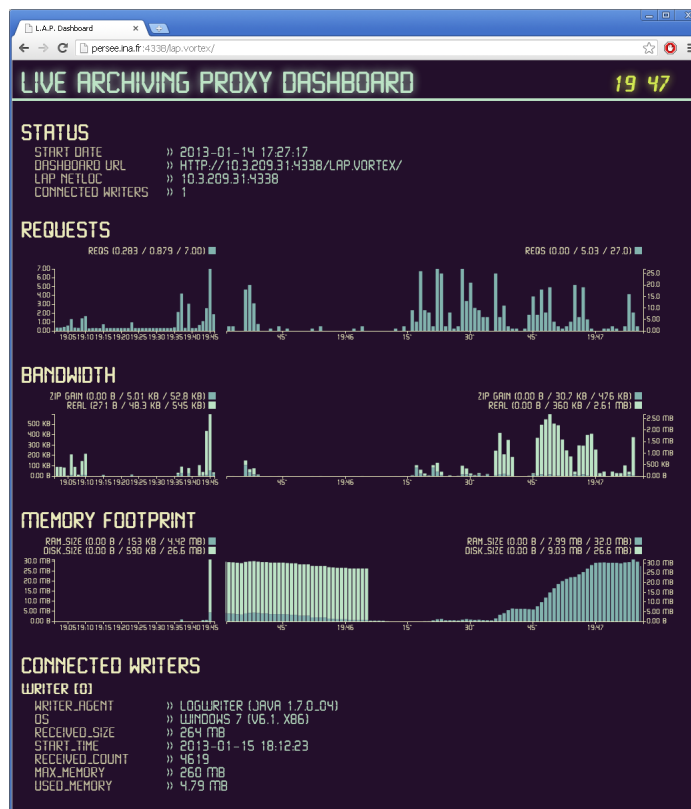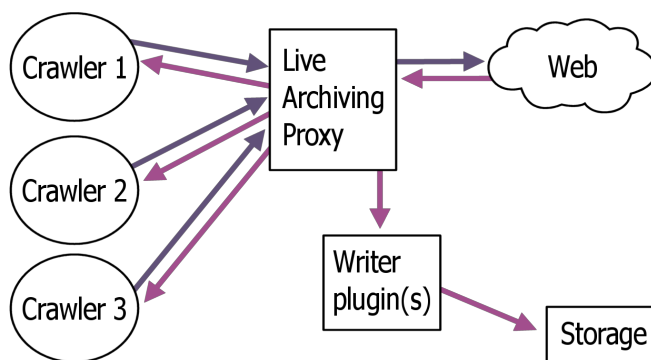# Live Archiving Proxy

*User Guide*

# General principles

The LAP (Live Archiving Proxy) is a single executable linux 64-bit binary file called "lap".
The LAP is a proxy through which you can access the Web like a normal proxy. All HTTP traffic that goes through the LAP is captured and sent to a dedicated writer. **If no writer is connected to the LAP, no archive will be created** and the used memory by the LAP will increase.



You can use an existing writer plugin or develop your own. As of now, this is the list of existing LAP writer plugins:

- WARC writer : download.
- Print writer (will only print captured URLs to console, not archive anything to disk) : download.

# Starting up a LAP with a WARC writer

## 1) Launch the LAP
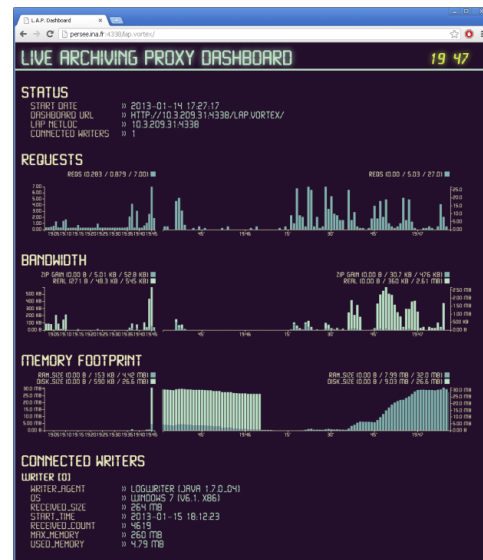
The LAP has two essential parameters:
- *WEB PORT* : the port the LAP listens on Web clients (default: 4338)
- *WRITER PORT* : the port the LAP listens on writers (default: 4365)

```
> lap --web-port WEB_PORT --writer-port WRITER_PORT
```





| Expected output | Dashboard interface |
|---|---|

Once the LAP is started, the dashboard is accessible with a browser at `http://LAP_HOST:WEB_PORT/`

## 2) Launch the WARC writer

The writer WARC has three essential parameters:
- *LAP HOST* : the IP or hostname of the LAP
- *WRITER PORT* : the port on which the LAP listens for writers
- *WARC DIR* : the directory in which WARC files must be created

```
> java -jar lap-writer-warc-1.0-SNAPSHOT-jar-with-dependencies.jar
LAP_HOST:WRITER_PORT --dir=WARC_DIR
```

The output should be similar to this :



## 3) Archive some stuff

Configure your crawler or any browser to use the proxy LAP_HOST on port WEB_PORT, and watch the LAP in action !

# Archiving HTTPS URLs

## 1) Pseudo-HTTPS mode

To make HTTPS traffic archivable, we make in non-encrypted between the Web client and the LAP. In order to do so, the Web client needs to alter outgoing HTTPS requests to make them HTTP, and append the fake ".HTTPS" TLD to the URL's host.
Here are a few examples of the alteration:

> https://example.com becomes
> http://example.com.https
>
> https://www.example.com/index.html becomes
> http://www.example.com.https/index.html

In this mode, the connection between the Web client and the LAP is unsecure and clear, but the connection between the LAP and the target Web server is true HTTPS.
The metadata sent to the writer will not contain the fake URL, but rather the actual HTTPS URL. From the writer's point of view, nothing really makes this request different except for the fact that the URL is HTTPS.

## 2) Altering URLs in PhantomJS to use pseudo-HTTPS mode

The following code will enable PhantomJS to alter HTTPS URLs on the fly to make them archivable by the LAP.

```
var page = require('webpage').create();

// LAP: alter HTTPS urls with fake .https TLD
var urlRegexp = new RegExp('^(https?)://([^/]+)(/.?)?$')
page.onResourceRequested = function(requestData, networkRequest) {
    var match = urlRegexp.exec(requestData.url);
    if (match.length >= 3 && match[1] === 'https') {
        networkRequest.changeUrl(
        'http://' + match[2] + '.https' + (match[3] === undefined ? '' : match[3])
        );
    }
};

var target = 'https://www.example.org/index.html';
page.open(target, function (status) {
    console.log('got status "' + status + '" for url "' + target + '"');
});
```

# LAP advanced configuration

**proxy**        The network location (host:port) of the HTTP proxy to be used by the LAP to access the Web. Using a caching proxy is useful for external bandwidth reduction.

**web-port**        The port on which the LAP listens for Web clients (default: 4338).

**writer-port**        The port on which the LAP listens for writer connections (default: 4365).

**digest**        The name of a digest (MD5, SHA1, SHA256). This digest will be computed by the LAP and sent to the writer. Use `metadata.getInfo("digest")` to read the digest in the writer.

**temp-dir**        The temporary directory used by the LAP to store big contents (over 5MB) before sending them to a writer (default: /tmp/LAP).

**bloom-netloc**        The network location (host:port) of a Bloom filter server used for deduplication.

# WARC writer advanced configuration

**dir**          Directory in which WARC file will be created

**prefix**       Prefix of the WARC files (default: LAP)

**deduplication**   Use deduplication mechanism (default: false).
An embedded database is used to avoid storing the contents multiple times by creating RT_IDX_REVISIT WARC records. The internal database is flushed every time the writer is started.

**compress**    Use WARC compression (default: false)

**max-file-size**  The maximum WARC file size in bytes (default: 1073741824 bytes = 1Gb). When the current WARC size reaches this size, it is closed and a new file is created.

**timeout**     Number of seconds before the connection to the LAP fails if the writer gets no response (default: 10).

**ispartof**    WARC Info : "ispartof" field (default: empty).

**description**  WARC Info : "description" field (default: empty).

**operator**    WARC Info : "operator" field (default: empty).

**httpheader**  WARC Info : "httpheader" field (default: empty).

**verbose**     Use verbose output (default: not enabled)