

# Automating Web History Analysis

---

**Michael Sonntag**

Institute for Information processing and  
microprocessor technology (FIM)

Johannes Kepler University Linz, Austria

[sonntag@fim.uni-linz.ac.at](mailto:sonntag@fim.uni-linz.ac.at)

# Why?

---

- Separate applications in the Internet are dying out
- Almost everything takes place (or is being moved) to the WWW
  - Exception: Apps on mobile devices
- Examples:
  - Social networks: "Native" web applications
  - Mail clients & calendars: Client program → Web-based interface
  - File storage: FTP → Sharehoster
  - Office software (documents, spreadsheets): Programs → Online services
  - Listening to music: Downloading → Streaming services (only interface via WWW!)
- Result: Web history gets more important, but also increasingly large & complex
  - Not a single browser window only, but multiple tabs + other applications

# Why “we”?

---

- Especially communication and collaboration has changed
  - Mandatory data retention: Don't use E-Mail → Web services!
  - Typical mode of communication for terrorists:
    - Draft E-Mail in web-based E-Mail system and store it, second person logs in and reads this same draft (and perhaps modifies it or generates another draft)
  - Typical exchange of illegal data:
    - Encrypt, assign random filename, upload to sharehoster
    - Pass on link in some way (perhaps: web forums)
    - Pass on password in some way (chat systems)
- All of these leave very few traces, and transport is encrypted (TLS!)
  - Inspect the local web history to find traces of this behaviour

# What?

---

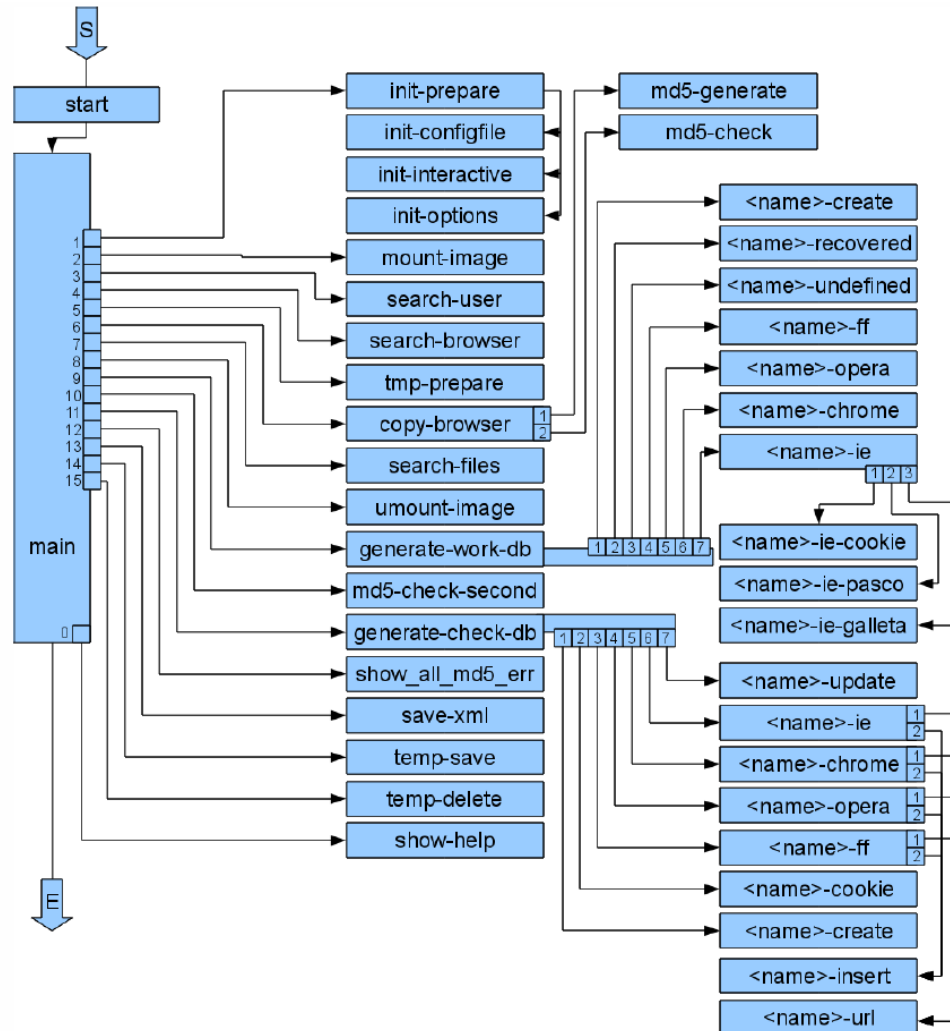
- Identifying the behaviour of users, based on their web activity traces
  - Traces remaining on the computer itself, i.e. no computer → no evidence!
    - Not: Data retention, online search!
- As much details as possible:
  - Which web pages were visited and what was their content?
  - Intentionality: Clicking or automatic JavaScript/PopUp?
  - Hints to external servers: Webmail, file storage etc.
  - Additional information: User names/passwords, visit count, ...
- From all browsers: A user might employ several, even on the same system
- Result: SW for automatically collecting and collating all this information
  - History, cache, and cookies → Timeline + sessions

# How?

---

- Command line tool: Important for automating the process
- Results: Text, CSV, XML, SQLite-DB → For extraction and evaluation
- Basic input: (Forensic – hopefully!) disk image
  - Not a copy of the files, but of all sectors of the disk → Requires mounting of file system
- Restriction to selected users possible if desired (useful for large servers)
- Reconstruction of deleted files optionally possible through external tool
- Validation of cache (+help for collation) against current web content (takes long!)
- Hash values are created & verified to protect against modifications
- Special handling of each type of browser → Simple extensibility
  - Supported: IE (requires two open source tools), Firefox, Chrome, Opera

# How?



# Show me!

- Exemplary output (XML version) of a single HTML file in the cache:

<pre> &lt;entry&gt;   &lt;id&gt;43&lt;/id&gt;   &lt;systemname&gt;windows_new&lt;/systemname&gt;   &lt;username&gt;MSCWin7&lt;/username&gt;   &lt;browsername&gt;ie&lt;/browsername&gt;   &lt;type&gt;http&lt;/type&gt;   &lt;url&gt;http://p8.friendscout24.de/e5e892b8a6eeb710c45b168bc5f2bd2a/0/78/57/29/14/74C3F09B5B4194909AD3FDBE820E142s.jpg&lt;/url&gt;   &lt;title&gt;://p8.friendscout24.de/e5e892b8a6eeb710c45b168bc5f2bd2a/0/78/57/29/14/74C3F09B5B4194909AD3FDBE820E142s.jpg   09/12/2009 19:24:09 09/12/2010 19:26:20   74C3F09B5B4194909AD3FDBE820E142s[1].jpg MSKM6NQ5 HTTP/1.0 200 OK Pragma: cached   Content-Type: image/jpeg Content-Length: 2123 X-Cache: HIT from p.friendscout24.de ~U:mscwin7&lt;/title&gt;    &lt;from_visit&gt;none&lt;/from_visit&gt;   &lt;last_visit_time&gt;1284312380&lt;/last_visit_time&gt;   &lt;visit_count&gt;0&lt;/visit_count&gt;   &lt;typed_count&gt;0&lt;/typed_count&gt;   &lt;hidden&gt;0&lt;/hidden&gt;   &lt;test_session_number&gt;1&lt;/test_session_number&gt;   &lt;test_session_time_start&gt;1284312381&lt;/test_session_time_start&gt;   &lt;test_session_time_stop&gt;1284312783&lt;/test_session_time_stop&gt;   &lt;test_cache_exist&gt;1&lt;/test_cache_exist&gt;   &lt;test_cache_same&gt;0&lt;/test_cache_same&gt;   &lt;test_url_exist&gt;1&lt;/test_url_exist&gt;   &lt;test_url_title_same&gt;0&lt;/test_url_title_same&gt;   &lt;test_cookies_set&gt;0&lt;/test_cookies_set&gt;   &lt;test_cookies_exist&gt;0&lt;/test_cookies_exist&gt;   &lt;test_cookies_time_ok&gt;0&lt;/test_cookies_time_ok&gt; &lt;/entry&gt; </pre>	<p>Internal number</p> <p>From which computer the data was collected</p> <p>Username of the inspected user</p> <p>Which browser</p> <p>http (=cache)/https (=cache)/cookie/visited(=history)/...</p> <p>The URL of this entry</p> <p>MSKM6NQ5 HTTP/1.0 200 OK Pragma: cached</p> <p>HIT from p.friendscout24.de ~U:mscwin7</p> <p>The stored HTTP response header</p> <p>ID of the previous URL</p> <p>The time of the last visit of this URL</p> <p>How often visited (0 = Not available)</p> <p>How often manually typed (0 = Clicked on link/Bookmark/...)</p> <p>Loaded in background (0 = not preloaded)</p> <p>Session number of reachability testing</p> <p>Start/End of testing session</p> <p>File has been found in the cache (here obvious!)</p> <p>File in cache differs from the one loaded in testing</p> <p>URL was found in testing (some page exists there)</p> <p>Title of test page download differs</p> <p>Test session did not set any cookies</p> <p>No cookies from test session found in cookies folder</p> <p>Number of cookies where the timestamp is similar to the ones downloaded during testing</p>
---	---

# What next?

---

- Adding further browsers, e.g. Safari, and especially those of mobile devices
  - Typically Webkit-based, so similar to Chrome/Firefox and each other
- Extracting more information: Difficult, as available data varies significantly between the browsers
- Automating the integration of several runs of the software
  - Several devices/computers of a single person
- Support for interpretation: Visualization, statistics, comparison to lists of known good/bad data (e.g. URLs or page content), search functionality (full-text index)
- Preprocessing: Assigning cookies to pages (note: need not be same domain → cache parsing!), web archive integration to complete missing cache elements, ...



# Conclusions

---

- Helpful program getting rid of very tedious and uninteresting work
  - Especially dealing with images and hash values: Necessary for forensic work!
- But still a lot of things to do: Interpretation
  - What does this mean? What did the person do? Anything in there suspicious?
    - Looking whether a specific URL was visited is trivial, but often the behavior is not nearly as easy to detect or conclusive!
  - Integration with other tools is still an effort, but at least it becomes possible
    - Existing tools often have GUI-output only!
- Sufficient for interpreting "web browsing"; doesn't work for web services
  - No "page" concept exists, not stored in "browser data"
    - Requires inspecting the web traffic itself (wiretapping)

# Thank you for your attention!

---

**Michael Sonntag**

Institute for Information processing and  
microprocessor technology (FIM)

Johannes Kepler University Linz, Austria

[sonntag@fim.uni-linz.ac.at](mailto:sonntag@fim.uni-linz.ac.at)