

EMM OSINT Suite – Improved Entity Recognition

Description

The EMM OSINT Suite is a desktop software package which consists of various tools based on JRC's research in open source text analysis and mining. The latest version features an improved entity recognition (extraction) module.

The software consists of the following core modules:

Data Acquisition

- Search – a component to extract search results from online search engines
- Crawler – a HTTP crawler module to harvest data from targeted web sites (“crawling”)
- Grabber – a HTTP client module to download text based or binary documents from web sites for further processing

Data Processing

- Text Extraction – extracts texts from different text based and binary formats (XML, TXT, PDF, MS Word, MS Excel, MS PowerPoint, Open Office)
- Entity Extraction – a set of modules to extract named entities from raw text. Entity types are people, organisations, locations, address information, VAT numbers and user defined custom types
- Category Matching – categorises text according to key word based category definitions

Data Analysis

- Reporting – a component to create reports for end users of for further external processing of extraction results
- Local Search – a local search index to provide full text search of downloaded artefacts
- Entity Browser – an analysis component to aggregate found entity data and allows browsing through the results.

The different tools are made available with a graphical user interface based on the Eclipse Rich Client Platform which is an open source toolkit for desktop applications.

Release Notes

The 2015 release contains the following improvements and bug fixes:

Improvements

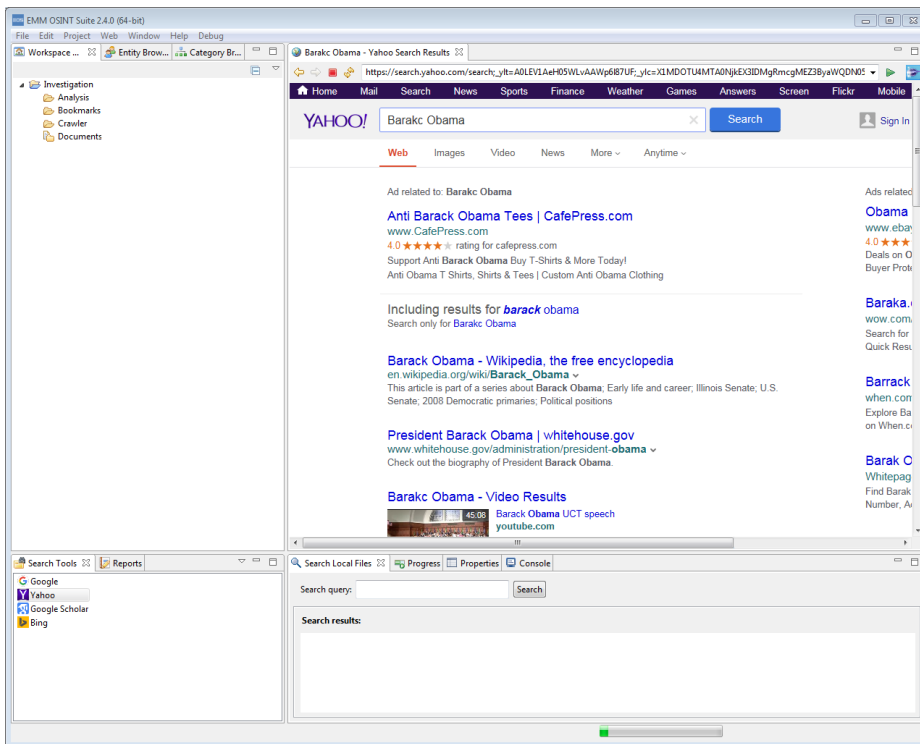
- **Improved Entity Recognition (Extraction)**
 - **Latest language resources and guessing patterns**
 - **New graphical view to test regular expression patterns (bot for BRICS and JAVA dialect)**
 - **New default custom pattern (for example Dutch zip code)**
- Improved Category Matching module
 - Find documents which match a specific keyword pattern
 - Updated domain specific language
- Internal analysis data model
 - Complete revision with substantial performance improvements and much less memory consumption
- Document handling
 - Improved document repository to enable on-demand translation service
 - Experimental import of complex document formats (for example MS Outlook PST)
- Automatic Software Update
 - In-place update from EMM server repository

Fixed issues

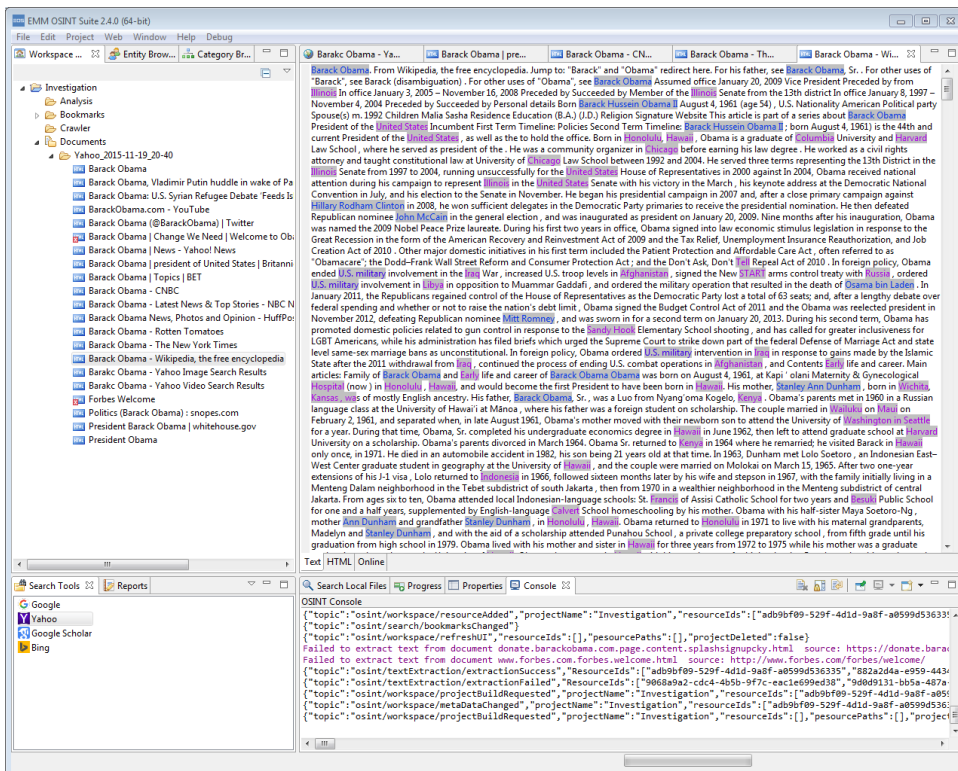
- Corrected graphical representation of entity relationships
- Updated Java Runtime
- Performance Improvements
- Completely revised build process based on Apache Maven and Tycho
- Fixed search result extraction patterns

Screenshots

Data Acquisition



Entity Extraction



Custom Entity Extraction

The screenshot displays the EMM OSINT Suite 2.4.0 (64-bit) interface. The main window is titled "Barack Obama..." and shows a configuration for a custom entity extraction rule. The "Node" pane on the left shows a tree structure for the rule, including "expressions", "declaration", "type", "id", "description", "expression", "regex", "name", "description", "validate", "output", and "output". The "Content" pane on the right shows the configuration details for the selected node, including the regular expression "(NL-)?[1-9][0-9]{3}[A-Z]{2}" and the description "Postal Zip Code (Netherlands)".

The "Test Regular Expressions" pane on the right shows the "Regular Expression:" field with the text "(NL-)?[1-9][0-9]{3}[A-Z]{2}" and the "Match" button. Below this is the "Text to match:" field and the "Found Matches:" field.

The "OSINT Console" at the bottom shows the following output:

```
OSINT Console
{"topic": "osint/workspace/projectChanged", "projectName": "Config", "resourceId": "osint/workspace/activeConfigProject"}
{"topic": "osint/catcher/activeConfigProject"}
{"topic": "osint/catcher/categoriesNotFound"}
{"topic": "osint/search/linkExtractorDefinitionsChanged"}
{"topic": "osint/workspace/resourceAdded", "projectName": "Config", "resourceId": "osint/workspace/activeConfigProjectContentChanged", "projectName": "osint/workspace/refreshUI", "resourceIds": [], "resourcePaths": []}
{"topic": "osint/report/templatesChanged"}
{"topic": "osint/catcher/categoriesNotFound"}
{"topic": "osint/search/linkExtractorDefinitionsChanged"}
Failed to initialize custom expression catches. Please check your custom...
```