# Next Generation Data Fusion Open Source Intelligence (OSINT) System Based on MPEG7

Mark Pfeiffer, Marco Avila, Gerhard Backfried, Norbert Pfannerer, Juergen Riedler

firstname.lastname@sail-technology.com

Sail Labs Technology AG

Vienna, Austria, EU

*Abstract*— **We describe the Sail Labs Media Mining System which is capable of processing vast amounts of data typically gathered from open sources in unstructured form. The data are processed by a set of components and the output is produced in MPEG7 format. The origin and kind of input may be as diverse as a set of satellite receivers monitoring TV stations or textual input from web-pages or RSS-feeds. A sequence of processing steps analyzing the audio, video and textual content of the input is carried out. The resulting output is made available for search and retrieval, analysis and visualization on a next generation Media Mining Server. Access to the system is web-based; the system can serve as a search platform across open, closed or secured networks. Data may also be extracted and exported and thus be made available in airgap networks. The Media Mining System can be used as a tool for situational awareness, information sharing and risk assessment.**

*Index Terms*—**Information Systems, Multimedia Computing, Speech Processing**

## I. INTRODUCTION

In today's world, an ever-increasing amount of information is being produced by the second and put on the internet and broadcast by TV- and radio stations. The online content stored on web-pages grows massively and at a constantly accelerating rate and is estimated to already exceed $1.6 \times 10^{20}$ bytes [1] – an increasingly large portion of which is multimedia (audio and video) content. Private users as well as professional entities place content on servers throughout the world in a multitude of formats and qualities and the percentage of content made available by users increases rapidly due to social network sites, blogs and the easier accessibility to the internet for an increasing number of users. TV and radio stations broadcast around the clock in a multitude of languages. To tap into this constant flow of information and make the multi-media contents searchable and manageable on a large scale, Sail Labs is developing a framework and system which allows the flexible combination of a variety of components for analysis of the different kinds of data involved. Information and clues extracted from audio- as well as video-tracks of multimodal documents are generated and stored for further analysis. The input files come from a variety of sources, whether pure audio or video and in a range of qualities from clear broadcast quality to low quality audio.

The information resulting from different threads of processing can be combined and used for analysis, information retrieval and visualization or decision support systems. The complete system is capable of processing vast amounts of data, such as those typically gathered from open sources, in a 24/7 mode of operation.

## II. SYSTEM DESCRIPTION

The Sail Labs Media Mining System consists of a set of technologies wrapped into components and models, which can be combined into a single system for end-to-end deployment. Together with the components, toolkits are delivered to allow users to update and refine models. The majority of models are based on statistics created from samples of data to allow for fast development and deployment as well as to facilitate multi-language development as much as possible. All textual input and output of the system is converted to UTF8 encoded Unicode which further facilitates multilingual processing.

The data run through a series of processing steps, which include speaker-identification (SID), language-identification (LID), automatic speech-recognition (ASR), face- and object-recognition as well as named-entity- and topic-detection. Automatic and human translation interfaces as well as notification interfaces are provided.

The components for processing may be grouped together in different configurations to allow for flexibility in setting up different kinds of systems. E.g. the cleaning and tokenization components are used as pre-processing steps for toolkits as well as post-processing steps for feeders extracting data from the web. Furthermore, not all components have to be present

in all configurations; subsets of components may be used for particular purposes, such as indexing text from the web only.

The result of processing are documents in a Sail Labs proprietary XML format or MPEG7 [2]. In case of MPEG7, different strategies have been implemented to fuse the results originating from different threads of processing, e.g. from visual-processing and audio-processing.

The resulting XML-files are uploaded together with a compressed version of original media files onto the Media Mining Server and made available for search and retrieval. The access to the server is web-based and can serve as a platform across open, closed or secure networks. The data may also be extracted and exported to a series of other media in order to make it available in airgap networks.

The overall architecture of the Media Mining System is a server-client one and allows for deployment of the different components on multiple-computers and platforms. Several Feeders, Media Mining Indexers and Media Mining Servers may be combined to form a complete system which can be used as a tool for situational awareness, information gathering and -sharing and risk assessment. Fig. 1 provides an overview of the components of the Media Mining System.
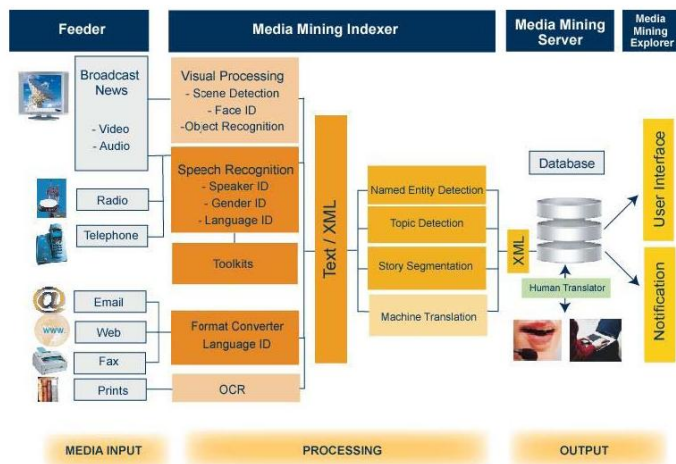


Fig.1. Components of the Sail Labs Media Mining System.

The following chapters provide an overview of the Media Mining System and describe in detail the individual components which make up the overall system.

### A. Feeders

The Feeders represent the interface of the Media Mining System to the outside world. For audio or mixed audio/video input a variety of formats can be ingested from external sources and processed by subsequent components. For example, the audio signal of a satellite receiver can be converted into the format required by the Media Mining Indexer. Likewise, to handle textual input, such as data

coming from Web-Pages or Emails, separate feeders exist which extract the data from these sources and pass them on to the text processing components.

The operation of the feeders can be controlled via an interface to allow for flexible arrangement in the flow of data.

For the processing of broadcast-news data and to facilitate a 24/7 mode of operation, a special feeder to link recording schedules with processing needs is provided.

### 1) Audio/Video Feeder

This feeder is based on the Direct Show framework and handles a variety of different input sources and formats. Files, sound-cards as well as TV- and satellite-cards are supported. Re-encoding is performed to provide Windows Media or Real Media output files which are uploaded to the Media Mining Server for storage and retrieval. The audio channel is passed on to the Media Mining Indexer for processing.

A special feature of this feeder is the built-in time-shifting capability which is used in connection with the Media Mining Indexer and translation technologies to provide sub-titling (with a minimal delay) of a TV program. This setup has been implemented and packaged in the award-winning product ROSIDS [3] (the Sail Labs Rapid Open Source Intelligence Deployment System).

### 2) Web-Collector

The Web-Collector is a special feeder to be used in the process of gathering data from internet sources such as web-pages or news-feeds (RSS-feeds). The Web-Collector provides a GUI to specify the task and monitor the operation of the collection process. Web-pages are scanned according to a variety of parameters which can be specified by the user. The final text extracted from internet sources is cleaned and tokenized before being handed on to the text-processing components.

### 3) Email-Collector

The Email-Collector is a special feeder used in the process of gathering data from email accounts. The Email-Collector provides a GUI to specify the tasks to be carried out, which user-accounts, servers, types of protocol etc is to be used. Emails are downloaded according to the description, cleaned and tokenized and finally passed on to the text-processing components.

### B. Media Mining Indexer (MMI)

The MMI forms the core of the Media Mining system. It contains a set of technologies, packaged as components, which perform a variety of analyses on the audio as well as the

textual content. Analysis of visual data can also be implemented and performed within the Indexer. However, Sail Labs does not provide visual technologies itself, but rather partners with research and commercial entities in this area. By providing a flexible, plug-in like architecture of the framework, we aim to make integration as easy as possible. Results of audio, textual and visual processing can be combined by enriching structures in an MPEG7 document. Models for processing a number of natural languages exist for the components of the MMI. For ASR currently more than a dozen models are available and this number is growing constantly.

### 1) Preprocessing of audio

After having been converted to the appropriate format by the feeder, the audio signal is processed and segmented for further analysis.

#### a) Signal Processing

The incoming audio is converted to a continuous sequence of audio-features, one per centi-second. Various normalization and conversion techniques are applied to the audio and features. Special processing can be added at this stage to handle specific formats such as processing of MP3 encoded audio, where feature-generation is carried out directly on the MP3 encoded data [4]. Finally, time-adjacent features are combined to yield the audio-features used by all subsequent components.

#### b) Segmentation

The incoming feature-stream is partitioned into homogeneous stretches or segments of audio. The segmentation stage uses models based on general sounds of language as well as non-language-sounds, such as breathing or lip-smacking noises, to determine the most appropriate segmentation point.

The content of a segment is analyzed with regard to the proportion of speech and non-speech contained in order to be able to decide how each segment should subsequently be processed. Segments classified as containing a sufficient amount of speech are passed on to the ASR component. Segments containing music can be processed by a component for commercial-detection or third-party components, such as music-analysis to determine the genre of a musical piece. A threshold can be set in order to prevent the system from running ASR on segments which contain exclusively music or at least a large proportion of music (typically jingles or commercials).

### 2) Speaker Identification (SID)

SID is applied to the segments produced by the segmentation step using a set of predefined target models. The set of target speakers typically comprises an initial set of persons of public interest. The set of target speakers can be extended or modified by users using a toolkit supplied by Sail Labs. In case a speaker's identity cannot be determined, the SID system tries to identify the speaker's gender. Data of the same speaker are clustered together and labeled with a unique ID, regardless of whether the speaker has been identified. Subsequent segments uttered by the same speaker are grouped together and marked as a *speaker-turn*. Data collected this way can then be used at a later stage to create new speaker-models.

### 3) Automatic Speech Recognition (ASR)

The Sail Labs speech recognition engine is designed for very large-vocabulary, speaker-independent, multi-lingual, real-time decoding of continuous speech. Recognition is performed in a multi-pass manner, each phase employing more elaborate and finer-grained models until the final recognition result is produced. Intermediate results can be made available if required. The recognition result is passed to the text-based components for further processing before being output to a proprietary XML document. The recognizer employs a time-synchronous, multi-stage search using Gaussian tied mixture-models, context dependent models of phonemes and word- and well as sub-word based n-gram models. The engine per se is language independent and can be run with a variety of models created for different choices of languages and bandwidths. Recognition is performed in a pipelined manner so as to guarantee maximum throughput and real-time behavior. Each recognition stage uses more elaborate and finer-grained models to refine intermediate results. Finally, text-normalization is applied to yield the sequence of decoded words or alternative sequences of decoded words (nbest-decoding), with time-tags and confidence-scores.

### 4) Language Identification (LID)

Within the scope of the Sail Labs Media Mining System, language identification is performed on audio- as well as on textural resources.
The audio language identification can be used to determine the language of an audio-document in order to allow processing of this file using a particular set of speech-recognition models, or, in a multi-server environment, route the audio-document to the appropriate server(s). Language dependent processing of speaker-turns in a multi-language conversation will be possible in a future version of the system.
Textual analysis of language is used to classify text before passing it on to the text-pre-processing components whose results are then passed on to the named-entity or topic-detection components. Likewise, textual analysis of language is applied in order to classify input for the language model toolkit.

## 5) Text-based technologies

The text-based technologies perform their processing either on the output of the ASR-component or on data provided by the text-pre-processing (text normalization) components.

### a) Text normalization

Textual normalization includes the pre-processing, cleaning, normalization and tokenization tools. Special handling of numbers, compound-words, abbreviations, acronyms, textual segmentation and normalization of spellings are all carried out by these components. Language dependent processing (sometimes even source-dependent processing) is performed. All textual data are processed as UTF8 encoded Unicode to facilitate multi-lingual development.

### b) Named entity detection (NED)

Detection of named-entities, such as persons, organizations and locations as well as numbers (telephone-numbers, currency amounts), is performed on the output of the ASR component, or, alternatively, on text provided by the text normalization components. The NED-system is based on patterns as well as statistical models over n-grams of words and is run in multiple stages. Sequences of words are searched and tagged according to the words, features defined over the words and the context the words appear in.

### c) Topic Detection (TD)

The topic-detection component performs two tasks: sections of text are classified according to a specific hierarchy of topics, and coherent stories are found by grouping together similar sections. The speaker-turns produced by the audio-segmentation stage form the basis for initial classification. Alternatively, paragraphs as produced by text normalization components can serve as the units for the initial classification stage. In a second phase, the already classified sections are compared to each other and adjacent sections whose content was classified as being similar are merged. The models used for TD and story segmentation are based on support vector machines (SVM) with linear kernels built on words. The hierarchy employed is one derived of the one used by Reuters and can be used across languages.

## 6) Text Indexer

The Sail Labs Text Indexer uses a subset of the components in the Media Mining framework to process different kinds of textual data. It does so by first using the appropriate feeder for access to the data source. The data are then normalized (using the text normalization components), its language identified and it is passed on to the text-technologies. Named entity detection, topic detection and story segmentation are performed and the output is provided in the same internal XML format as the output produced by ASR.

## 7) Toolkits

The majority of models used in the Media Mining System is of a statistical nature. As customers might possess sufficient amounts of data to augment, extend or replace these models, we have been putting special emphasis on the development of toolkits to accompany the components and models which form the Media Mining System. While ideally we would like to allow users to re-train all the models involved, there are currently two toolkits available to allow user intervention. However, in cooperation with customers special models can be created by Sail Labs.

### a) Language Model Tookit (LMT)

The Language Model Toolkit (LMT) allows users to refine or extend the scope of the ASR-component of the Media Mining System. The LMT provides a GUI which allows users to build language models using their proprietary text resources. The ASR-component's vocabulary and language model can be extended by adding new words and contexts to them. The vocabulary and pronunciations can be adjusted according to particular requirements such as geographically or socially motivated uses of the system, e.g. by providing pronunciations for a certain geographical area or including slang- or swear-words typically used by the target group of speakers. A host of parameters can be set to best fit the model to the target domain. The LMT builds the model and installs it, so the model can be used immediately after creation. All tasks performed via the GUI are also available through a simple API which allows for easy embedding of the LMT into an existing environment and work-flow.

### b) Speaker Identification Toolkit (SIT)

The Speaker Identification Toolkit allows users to refine or extend the scope of the SID component of the Media Mining System. Data collected for speakers during operation of the Media Mining System can serve as the basis to train new speaker models. By providing feedback about a previously unknown speaker's identity, the data collected during processing of this speaker's audio can be used to create a new speaker model for them. A few minutes of audio per speaker usually form a good basis for the creation of such a model. The SID system can be trained in an incremental manner, so future addition of persons (e.g. caused by a change of government) can be accommodated easily. Trained models can be added to the existing set of speakers or existing speakers be replaced with the newly created models.

## 8) Machine Translation

Sail Labs is not active in the area of automatic translation and rather partners with companies active in the field, such as Sakhr [5] among others.

### 9) CAVA Framework

Sail Labs has been participating in a series of Austrian Research Projects [6] with the aim of creating technologies to allow for long-term, unsupervised adaptation of a system deployed in the real-world. Clearly, the models used for each of the sub-systems described above age over time. Some of them age more rapidly than others, but all models suffer from the fact that they run out-of-date at some point in time. For example, ASR vocabularies change over time. Persons and words come into and go out of usage as events develop over time. Within the CAVA framework (**C**ontinuous **A**utomatic **V**ocabulary and Language Model **A**daptation) existing components are linked together and additional components are created to allow for a flexible and autonomous mode of operation and adaptation of systems deployed at customer-sites. Information extracted and gathered using the text feeders, normalization- and text-processing tools, together with the base-models provided by Sail Labs, is used to update language models and vocabularies automatically. The point in time at which such an update seems most appropriate as well as the manner of how to exactly perform the update is determined by the system. A sudden change in the vocabulary of selected sources may provide a signal to adjust the vocabulary and rebuild speech recognition models.

Another source of data to include is formed by corrections made by users via the Media Mining Explorer. Corrections are gathered and kept track of in order to benefit from them and to trigger rebuilding steps. Depending on the amount and recency of corrections, words may become part of the active vocabulary.

### C. Media Mining Server (MMS)

The Media Mining Server comprises the actual server, used for storage of XML and media files, as well as a set of tools and interfaces used to update and query the contents of the database.

### 1) Media Server

The actual server provides the storage for the XML index files, the audio and the video content. It makes use of a database which provides the basis for all search and retrieval functionality. Currently, this database is Oracle 10g [7], as this product provides all features needed for state-of-the-art and timely information retrieval. Complex queries as well as summarization are part of the suite of technologies offered.

### 2) Translation

Different types of and interfaces to translation facilities are offered by the MMS. Parallel translations, potentially in several languages, can be created for a transcript as it is uploaded to the server, via integration of automatic translation engines.

Searches are applied to the original transcript as well as to the translated versions of all documents. Terms mentioned in a search can be translated on a per-term basis, so that documents mentioning the equivalent term in a different language will also be retrieved. This simple term-based translation can also be applied to the results produced by a query.

An interface to human translators is supplied as well. Users of the system may trigger a request for human translation in cases where the simple term-based or automatic translation is not sufficient and the need arises to understand a specific document more thoroughly. The translated version is then re-integrated into the database automatically, i.e. transparently to the end-user.

### 3) User Interaction and Feedback

Sail Labs provides a set of tools to let users interact with and update the contents of the database.

Through a web-browser users can interact with the MMS, perform queries, download content, associate stories with users, request translations and add annotations to the information stored. Statistics about the occurrence of terms can be visualized and scanned.

Fig.2 depicts the query/result interface with a transcription of Al-Jazeera and an automatically generated English translation. Speaker turns are annotated to the right, the video is played back in sync with the audio in the upper-right corner.



Fig. 2. One result returned by a query (with translation)

### a) Queries

Queries can be performed using a combination of search terms targeted at particular fields or free text and logical combinations of such terms. They can be tailored to address only specific portions of the data stored on the server, such as limiting them to a certain channel, period of time, certain speakers, topics or locations.

Users can attach special keywords to documents stored in the database and later also perform searches specifically on these keywords.

Queries can be stored and later used for automatic notification of users to allow for rapid notification when new documents matching a particular profile appear on the server. The MMS thus provides functionality to associate users with content, i.e. documents which are typically of interest only to a particular user or groups of users. A mobile client has been developed for use with notifications within the scope of the Reveal This project [8].

The results of queries are presented according to the segmentation produced by the MMI. Information, such as the names of speakers, named-entities, topics associated with a document as well as information produced by the visual indexing components such as keyframes, names of persons whose faces were identified or objects recognized are displayed along with a transcript of the associated audio. In case of text-only documents, certain types of information will be missing. Playback of audio and video content can be triggered on a per-segment basis or for the complete document. The words are highlighted synchronously with the playback of the audio/video.

*b)*        *Other types of user interaction*

When viewing results of queries, it is inevitable to come across errors, e.g. in the transcripts produced by ASR. Users can correct the text and topic of documents or the names of speakers. These corrections are stored for potential use in the vocabulary-updates and will also be reflected in the database.

Furthermore, users can upload multi-media content. Processing of these files will be triggered automatically and the results enter into the usual flow of data.

Media files and transcriptions can also be downloaded by users to prepare special reports/documents/dossiers for offline analysis.

## III.   TYPICAL INSTALLATION SCENARIO

Typical installations of the Media Mining System include a set of Feeders, Media Mining Indexers and Media Mining Servers. The exact number of each of these components varies with the amount of data to be processed. In typical 24/7 scenarios, a set of audio/video feeders is connected to satellite, TV or radio input. The feeders provide a constant flow of input to a set of Media Mining Indexers. These indexers process the incoming data stream and continuously send their outputs to the Media Mining Servers. Multiple servers may be used for fail-safety and to increase the bandwidth for access by many users. In addition, text-feeders may be used to provide input from Web in parallel to the multi-media input.

## IV.   CONCLUSION

The system presented is a current state-of-the-art end-to-end OSINT system and in the process of becoming a part of the state-of-the-art next generation in OSINT systems. Experience gained through projects with various governments in Europe, Africa, the Middle East, and various US- and Asian clients, who have shared their views and wishes on what the required features are have helped Sail Labs develop the Media Mining System and framework. It is necessary to have an open minded and also open interfaced solution. This has been one of the main requests by the government entities within the EUROSINT Forum Technology Gaps group [8]. Open Source Intelligence solutions must be interoperable with each other's solutions. It is simply not acceptable for governments or other entities to be locked into one proprietary non-interoperable system that stands alone similar to a silo. This is regarded as a waste of resources and simply a luxury that is not affordable given the high stakes that are involved. Consequently we aim to make our framework and systems as flexible as possible. By allowing plug-ins at various levels we provide the flexibility required by customers as well as allowing different technologies (possibly from different vendors) to be combined.

## REFERENCES

[1]   C.H.Best, "Open Source Intelligence", Joint Research Centre, European Commission, reference to data from IDC
[2]   http://en.wikipedia.org/wiki/MPEG-7
[3]   http://en.wikipedia.org/wiki/ROSIDS
[4]   http://www.ist-divas.eu
[5]   http://www.sakhr.com
[6]   http://www.coast.at
[7]   http://www.oracle.com
[8]   http://www.reveal-this.org
[9]   http://www.eurosint.eu (Vienna meeting 2007 at OSCE)