

# Ontea: Pattern based Semantic Annotation Platform



Michal Laclavík, Martin Šeleng, Ladislav Hluchý

laclavik.ui@savba.sk, <http://ontea.sourceforge.net/>, <http://ikt.ui.sav.sk/>

Institute of Informatics, Slovak Academy of Science, Dúbravská cesta 9, 845 07 Bratislava, Slovakia

## Motivation

To create semantic meta data from texts or documents

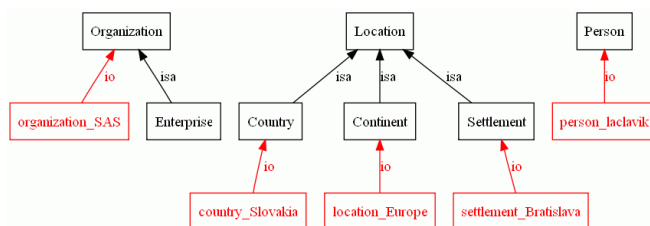
## Approach

Even unstructured text documents contain patterns, which can be used to extract various objects from text. Ontea is platform which exploits pattern approach to find or create key - value pairs from text. Pairs can be further processed and transformed to ontology instances or other metadata type.

## Pattern Types

Regular Expressions - current implementation  
Further possibilities:

- Integration of existing pattern solutions such as C-PANKOW or GATE
- Advanced patterns used in Information Extraction and Information Retrieval
- Keyword based
- XPath - for XML or HTML documents



Example	Text	Patterns - Regular Expressions
1 English	Michal Laclavik works for Slovak Academy of Sciences located in Bratislava, the capital of Slovakia	<code>\\b(\\p{Lu}[a-z]+ \\p{Lu}[a-z]+ \\p{Lu}[a-z]+)\\b</code>
2 English	Car maker KIA Motors decided to build new plant in Slovakia near town of Zilina. It is its first plant in Europe. Contact: Kia Motors Slovakia, Ltd. P.O.Box 2 01301 Teplicka nad Vahom Slovakia	<code>(in near) +(\\p{Lu}\\p{L}+) Location (city town) of (\\p{Lu}\\p{L}+ *\\p{Lu})\\p{L}*) Settlement \\b(\\p{Lu})[-&amp;\\p{L}]+[ ]*[-&amp;\\p{L}]* ]*[-&amp;\\p{L}]*[, ]+ (Inc Ltd)[.\\s]+ Organization</code>



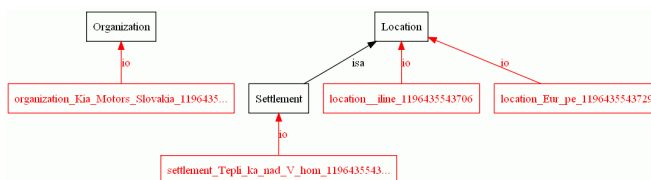
Example	Method	Annotation Results
1	Ontea	Michal Laclavik <i>Person</i> Slovak Academy <i>Organization</i> Bratislava <i>Settlement</i> Slovakia <i>Country</i>
2	Ontea create	Slovakia <i>Country</i> Zilina <i>Settlement</i> Europe <i>Continent</i> Kia Motors Slovakia <i>Organization</i>

## Used Technology

- Regular Expressions
- RDF, OWL
- Java

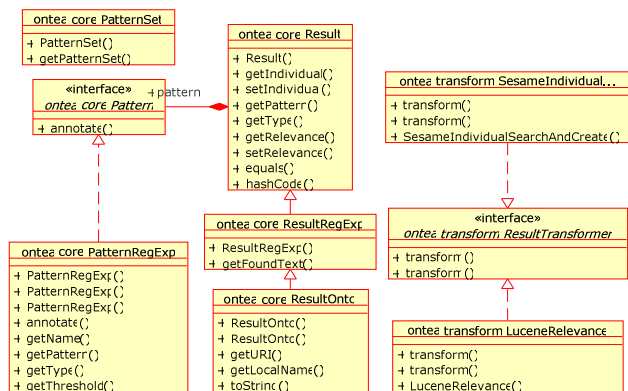
## Used Tools

- Lucene
- Hadoop
- Sesame, Jena



## Features

- Identification of concept instances from the ontology
- Automatic population of ontologies with instances
- Identifying relevance, when creating instances using information retrieval techniques
- Large scale semantic annotation of documents or texts using Google's MapReduce architecture.



## Architecture

- **ontea.core.Pattern:** Interface to adopt different pattern annotation techniques. Current implementations include regular expression pattern matching
- **ontea.core.Result:** Class representing results of pattern annotation - instances of defined type
- **ontea.transform.ResultTransformer:** Interface transforming results of annotation to different type or quality of results e.g. concrete ontology mapping, knowledge base implementation or result quality checking

## Integration with External Tools

- **Naliti:** Text Language Identification
- **Morphonary:** Lemmatization of Slovak developed at UPJŠ Košice
- **Lucene:** New instance relevance identification.
- **Hadoop:** Large scale semantic annotation using MapReduce Architecture
- **Sesame, Jena:** Transformation of found key - value pairs into RDFS or OWL instances in Sesame or Jena API