

# AnnoMarket – Multilingual Text Analytics at Scale on the Cloud

Marin Dimitrov<sup>1</sup>, Petar Kostov<sup>1</sup>, Hamish Cunningham<sup>2</sup>, Ian Roberts<sup>2</sup>, Philippe Rigaux<sup>3</sup>, Helen Lippell<sup>4</sup>

<sup>1</sup> Ontotext AD, Bulgaria

{marin.dimitrov, petar.kostov}@ontotext.com

<sup>2</sup> Department of Computer Science, University of Sheffield, UK

{h.cunningham, i.roberts}@dcs.shef.ac.uk

<sup>3</sup> Internet Memory Research SAS, France

philippe.rigaux@internetmemory.net

<sup>4</sup> The Press Association Ltd, UK

helen.lippell@pressassociation.com

**Abstract.** AnnoMarket is an FP7 funded project providing an open platform for cloud-based text analytics services and language resources acquisition. Providers of text analytics services and language resources can deploy and monetize their components via the AnnoMarket platform, while users can utilize such available resources in multiple languages and in various domains in an on-demand, pay-as-you-go manner.

**Keywords:** text mining, cloud computing, software-as-a-service, linked data

## 1 Introduction

AnnoMarket<sup>1</sup> is an FP7<sup>2</sup> project that aims to revolutionise the text analytics market, by delivering an open marketplace for pay-as-you-go, cloud-based text mining resources and services, in multiple languages. The current services available on the AnnoMarket marketplace<sup>3</sup> are applicable to a wide set of business cases, e.g. large-volume multi-lingual information management, business intelligence, social media monitoring, customer relations management.

The Software-as-a-Service delivery model adopted by AnnoMarket reduces the complexity of deployment, maintenance, customisation and sharing of text processing resources and services by SMEs and developers. The marketplace currently provides various services for multilingual information extraction and semantic annotation, sentiment detection, as well as multilingual web corpora, public language resources and LOD datasets.

---

<sup>1</sup> <https://annomarket.eu/>

<sup>2</sup> The AnnoMarket project is funded by the European Commission under the 7<sup>th</sup> Framework Programme, Project No. 296322

<sup>3</sup> <https://annomarket.com/>

## 2 AnnoMarket Platform

The AnnoMarket platform is comprised of components for language resource acquisition, scalable and elastic text mining over large volumes of data, usage monitoring and quota enforcement, as well as billing and online payments.

### 2.1 Language Resource Acquisition

The language resource acquisition component of the AnnoMarket platform is based on the large scale web crawling infrastructure providing user-defined, on-demand crawls at a large scale. In addition to corpora crawled on-demand from the web, the language resource acquisition component provides integration with the Common Crawl<sup>4</sup> dataset.

### 2.2 Multilingual Text Mining Services

Various multilingual text mining services are currently deployed and ready to use via the AnnoMarket platform. The current set of services includes more than 30 different text processing pipelines covering 17 languages.

The text processing pipelines vary from low level ones (stemmers, part-of-speech taggers, noun phrase chunkers and parsers), to general purpose pipelines (named entity recognisers) and domain specific pipelines (for the bio-medical domain, news publishing domain, or sentiment analysis over social media).

### 2.3 Marketplace

The marketplace provides an eShop<sup>5</sup> where customers can explore the catalogue of available text analytics services, language resources and datasets as well as additional processing resources available on-demand via the platform (e.g. an LOD server hosting Freebase, DBpedia and GeoNames; indexing servers, etc.). All products deployed on the marketplace provide information about their functionality and the associated usage and pricing terms.

### 2.4 Cloud Platform

The AnnoMarket platform is currently deployed on the Amazon Web Services<sup>6</sup> public cloud and it utilizes various cloud services for storage, computing, scalability, and a design for a multi-datacenter deployment for improved availability.

---

<sup>4</sup> <http://commoncrawl.org/>

<sup>5</sup> <https://annomarket.com/shopfront>

<sup>6</sup> <https://aws.amazon.com/>