# AIDA: combining text mining and e-science for bioinformatics

Marco Roos[1,a], Scott Marshall[1,b], Sophia Katrenko[1], Edgar Meij[1], Willem van Hage[2], Frans Verster[1], Pieter Adriaans[1]

Adaptive Information Disclosure: 1) Institute of Informatics, Faculty of Science, University of Amsterdam, The Netherlands;
2) TNO Industrie & Techniek / Vrije Universiteit Amsterdam, The Netherlands
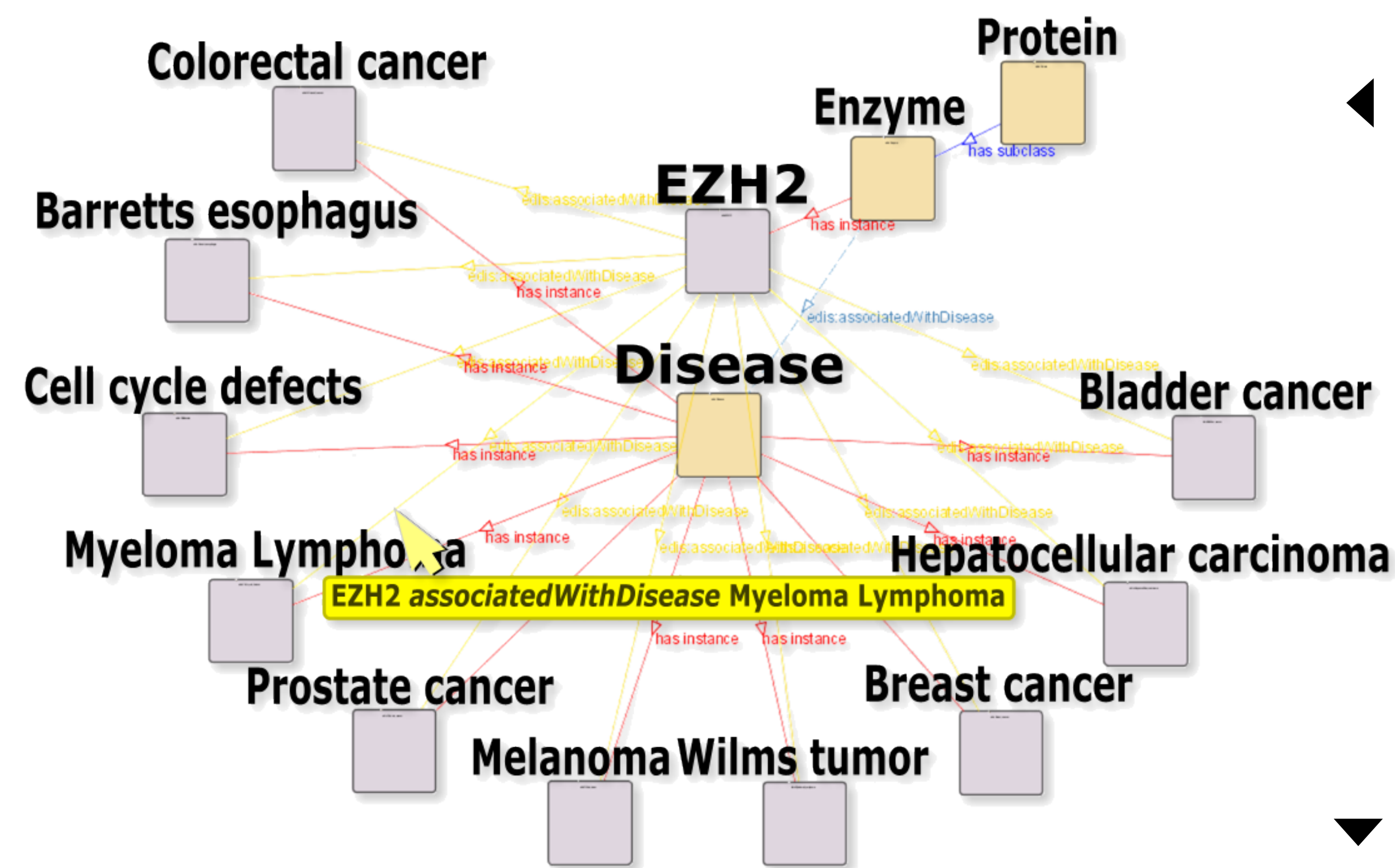a) Project or Area Liaison (PAL) for OMII-UK (e.g. Taverna)    b) Participating in W3C Semantic Web Health Care and Life Sciences Interest Group ( http://www.w3.org/2001/sw/hcls/ )

Correspondence: marshall@science.uva.nl, Website: http://adaptivedisclosure.org

## Introduction

### Our objective is to provide heuristic support for hypothesis construction from literature

We start with a seed of knowledge, a 'proto-ontology', that we want to extend. For example, information from a review article about histones and disease.



◁ Snapshot of Jambalaya plug-in for Protégé/OWL showing relationships between Enhancer of Zeste Homologue 2 and various diseases.
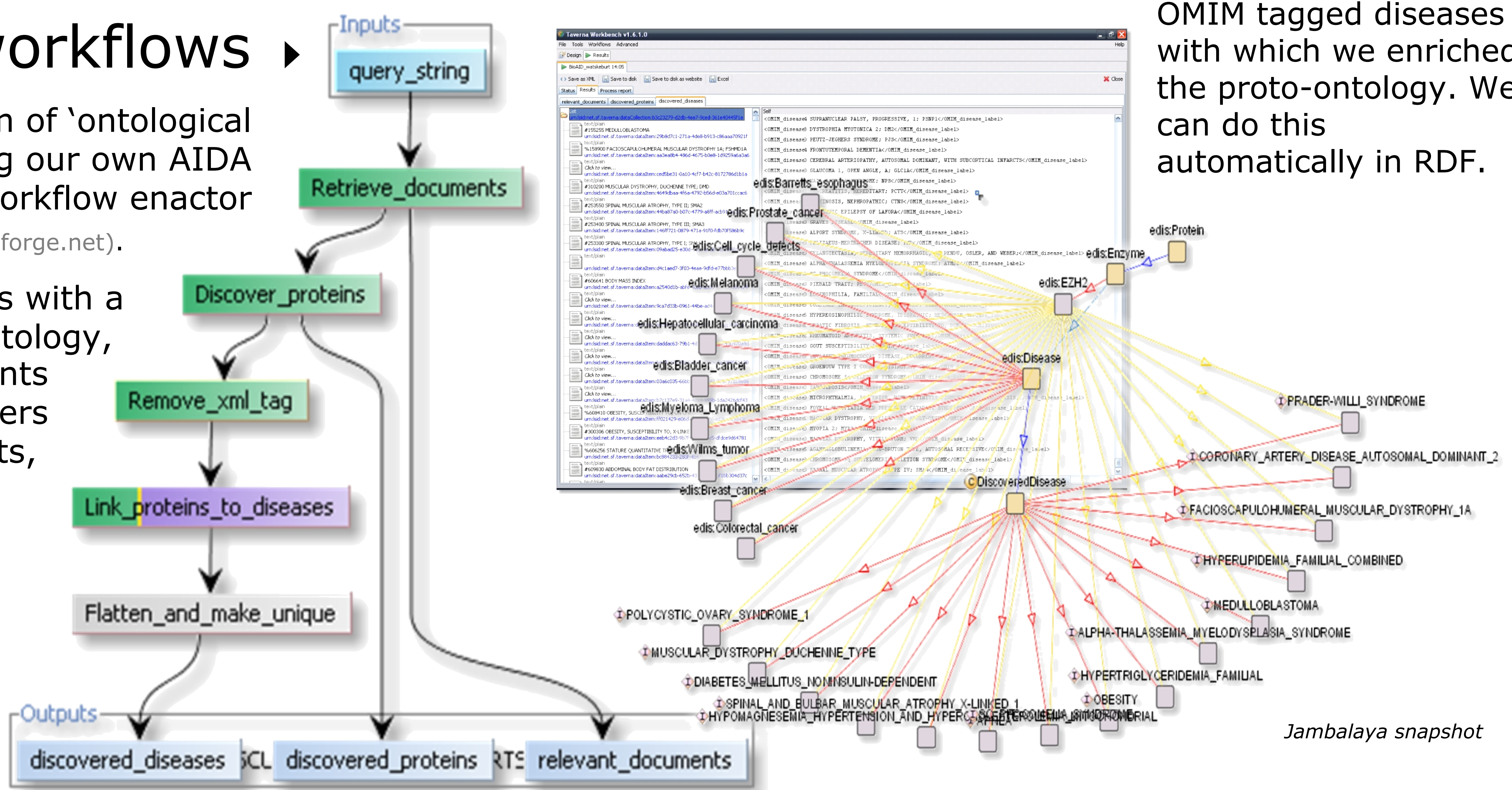From table by Moss TJ, Wallrath LL, Mutation Research, Vol. 618, pp. 163-174.

## Text mining workflows ▶

We implemented a form of 'ontological enrichment' by connecting our own AIDA services and others in the workflow enactor tool Taverna (http://taverna.sourceforge.net).

The workflow here starts with a query from the proto-ontology, retrieves relevant documents from medline, discovers proteins in the abstracts, and then links these proteins to diseases using the OMIM database

OMIM service from the Center for Information Biology and DNA Data Bank of Japan, National Institute of Genetics (see http://xml.nig.ac.jp)
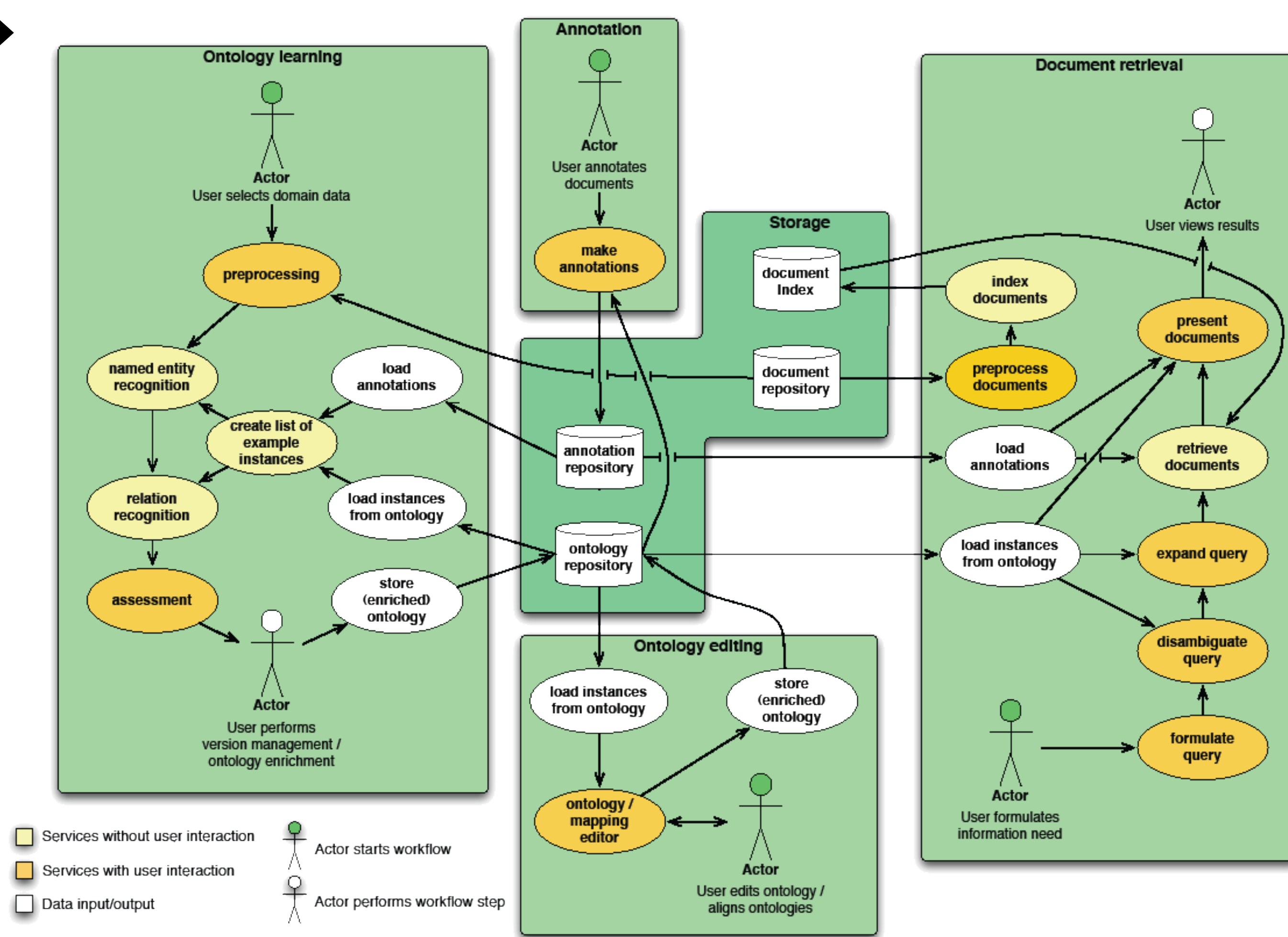


## Results

▼ The workflow produces OMIM tagged diseases with which we enriched the proto-ontology. We can do this automatically in RDF.

*Jambalaya snapshot*

## The AIDA toolbox ▶

An alternative to building special purpose bioinformatics tools is to use *web services* within a general computing environment such as the workflow tool *Taverna*.

Our workflows apply this 'e-science' approach to text mining. The components come from our 'AIDA* toolbox', which contains services that can be flexibly combined for a variety of applications, including various forms of text mining.

*Adaptive Information Disclosure Application



## Discussion

Our example text mining workflow supports hypothesis construction by enriching an ontology with hypothetically related elements that were discovered from literature. We can explore various text mining strategies by adapting workflows and by adding new services from the AIDA toolbox, such as machine learning services to make discoveries in terms of one's own ontology, and semantic web services to steer data handling in terms of ontologies. The proportion of true positives returned (recall) can be improved, for instance, by incorporating synonym services. The AIDA toolbox is of general use. For instance it is also used in food informatics applications.

## Acknowledgements and availability

nbic