



ANT The System Architecture of an Academic Search Engine

José Devezas

joseluisdevezas@gmail.com

INESC TEC & FEUP InfoLab

MAP-i 2016/2017



Universidade do Minho



universidade
de aveiro





Contents

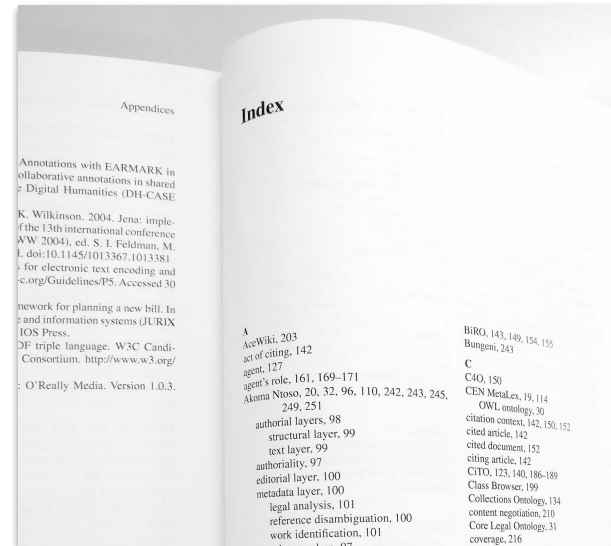
- Introduction
 - Search using keyword-based matching
 - Search using semantic matching
 - Requirements of an EOS engine
- ANT
 - An academic EOS engine
 - ANT: Ad hoc search of eNtities and Text
 - Search engine architecture
- Conclusions
 - Final remarks
 - Related projects

Introduction

What is entity-oriented search?

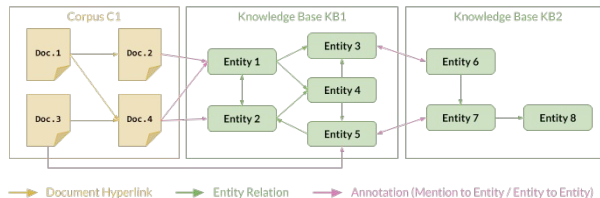
Search using keyword-based matching

- Modeled after the back-of-the-book index.
- Finding relevant content involves:
 1. Selecting one or several keywords;
 2. Jumping to the indicated pages;
 3. Reading passages and using knowledge, either internal or external to the book, to assess the relevance.



Search using semantic matching

- Closer to the user's information need.
- Requires interpretation of query meaning and document semantics.
- And the combination of unstructured and structured data from corpora and knowledge bases.

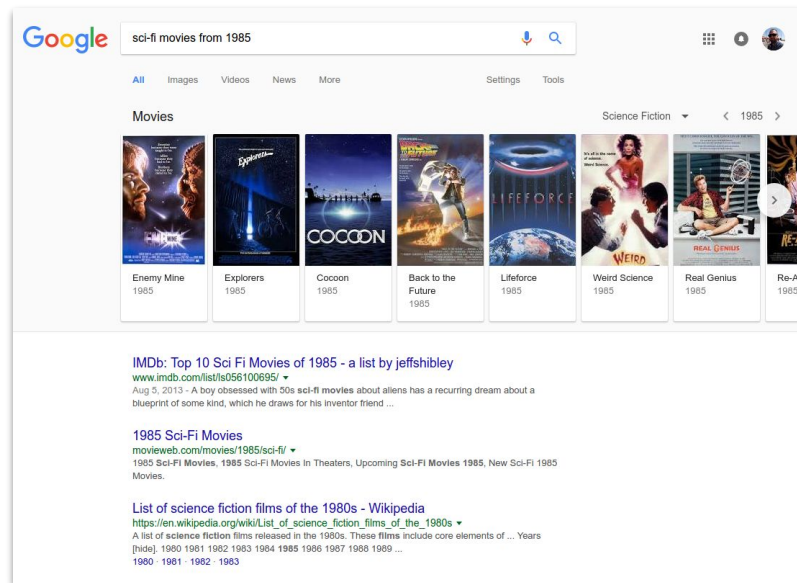


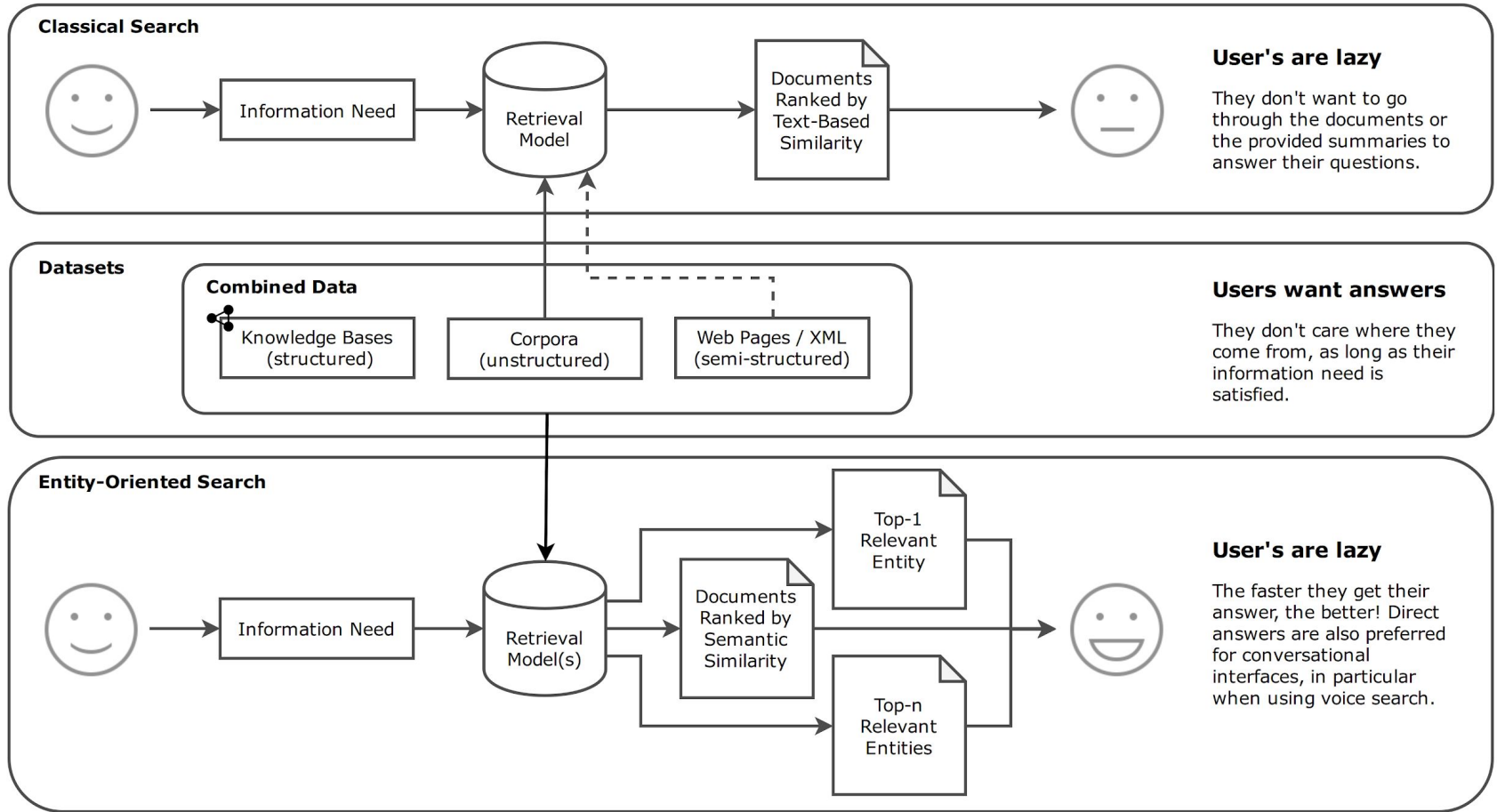
The screenshot shows a Google search result for the query "president of portugal 2018". The search results page displays the following information:

- Search Query:** president of portugal 2018
- Results:** About 187,000,000 results (0.60 seconds)
- Result Title:** President of Portugal (2018)
- Result Image:** A portrait of Marcelo Rebelo de Sousa.
- Result Text:** Marcelo Rebelo de Sousa
- People also search for:** A list of seven other presidents of Portugal: António Costa, Eduardo Ferro Rodrigues, Aníbal Cavaco Silva, José Sócrates, Pedro Passos Coelho, Rui Rio, and Jorge Sampaio.
- More about Marcelo Rebelo de Sousa:** A section with a "Claim this knowledge panel" link and a "Feedback" link.
- President of Portugal - Wikipedia:** A link to the Wikipedia page for the President of Portugal.
- Wikipedia Snippet:** The current President of Portugal is **Marcelo Rebelo de Sousa**, who took office on 9 March 2016.
 - Term length:** Five years; Renewable once, con...
 - Salary:** €93,364.74 (2015); (€6,668.91/month)
 - Residence:** Belém Palace
 - First holder:** Manuel de Arriaga
 - Role:** Powers · Election · 2016 presidential election

Search using semantic matching

- It becomes possible to, more adequately, answer queries like:
 - [president of portugal 2018]
 - [sci-fi movies from 1985]
- In addition to ranking documents containing the keywords;
- An entity or list of entities is directly provided as the answer.







Requirements of an EOS engine

- Integrate text, entities and their relations.
 - Where to integrate?
 - Index level;
 - Ranking level.
 - What are our technological options?
 - Inverted index;
e.g., Apache Lucene.
 - Triple store;
e.g., OpenLink Virtuoso.
 - Hybrid technology?
Higher-risk: still being researched.
- Match query and documents/entities using all available information for ranking.
 - Query segmentation and semantic tagging;
NER in query / linking entities to the knowledge base.
 - Document annotation.
NER in documents / linking entities to the knowledge base.

ANT

Search engine architecture.





An academic EOS engine

Depending on the query, results can be:

- Documents*
 - Retrieved using semantic information (entities and their relations).
- Entities
 - A specific one, a list, or both.
 - Retrieved by name, type, or another description.
 - Representing attributes or relations.

Entities can be:

- Students
- Staff
- Departments
- Rooms
- Curricular Units
- Courses
- News*



Ad hoc search of eNtities and Text.

- Designed to support the five query categories, as defined by Pound et al. (2010):
 - Entity;
 - Type;
 - Attribute;
 - Relation;
 - Keyword.
- Based on two Lucene indexes:
 - Query analysis index;
 - Entity index.
- And a Virtuoso quad store:
 - Useful for attribute and relation queries.

[Todos](#) [Notícias](#) [Estudantes](#) [Salas](#) [Pessoal](#) [Cadeiras](#) [Cursos](#) [Departamentos](#) | [Ferramentas de Pesquisa](#)

26340 resultados (16.51 segundos) 



José Luís da Silva Devezas

Estudante https://sigarra.up.pt/flup/pt/vld_entidades_geral.entidade_pa...

Faculdade de Belas Artes da Universidade do Porto (FBAUP) (mais 4)

Curso: **Doutoramento em Informática**

Código: 200303288





Vítor Bruno dos Santos Devezas

Estudante https://sigarra.up.pt/icbas/pt/vld_entidades_geral.entidade_p...

Instituto de Ciências Biomédicas Abel Salazar (ICBAS)

Curso: **Mestrado Integrado em Medicina**

Código: 200706081





Diana Maria Rodrigues Alves Devezas

Estudante https://sigarra.up.pt/flup/pt/vld_entidades_geral.entidade_pa...

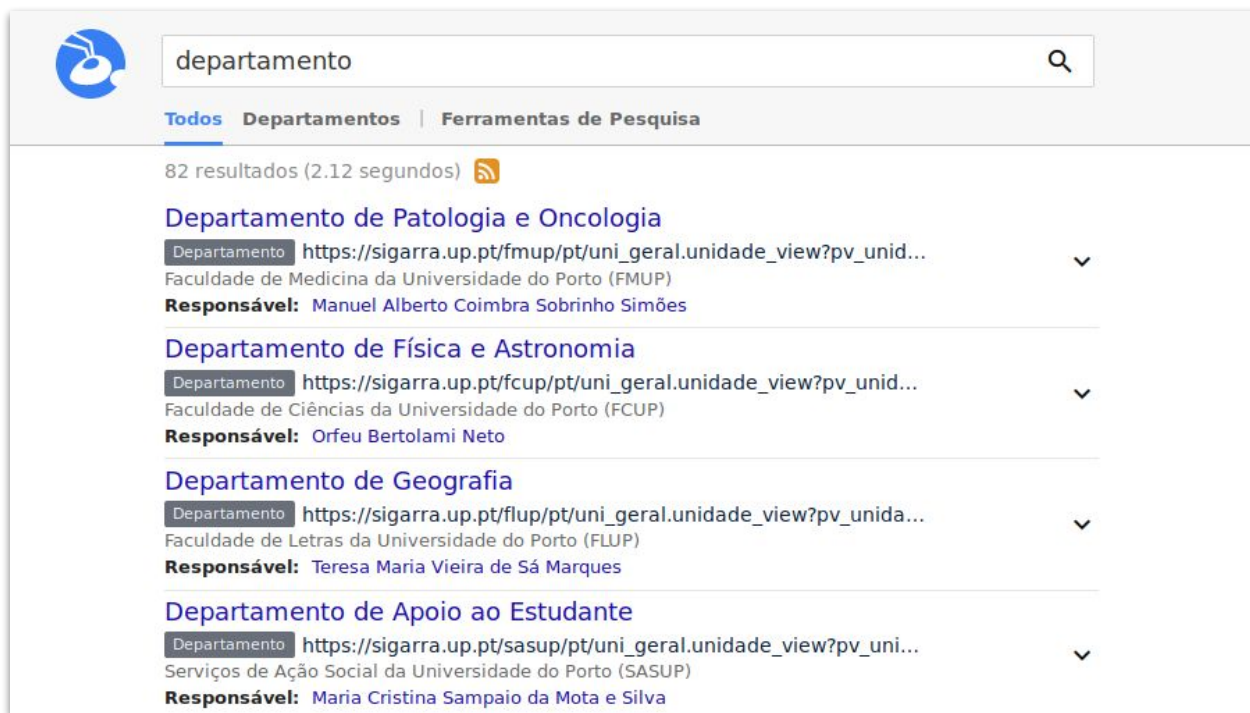
Faculdade de Letras da Universidade do Porto (FLUP) (mais 1)

Curso: **Direito**

Código: 200300607



Entity query. *The intention of the query is to find an entity.*



The screenshot shows a search interface with a logo on the top left, a search bar containing the word "departamento", and a magnifying glass icon. Below the search bar are three tabs: "Todos", "Departamentos" (which is selected), and "Ferramentas de Pesquisa". The results section shows "82 resultados (2.12 segundos)" with an RSS icon. Four department entries are listed, each with a title, a URL, a description, and a responsible person. Each entry has a downward arrow icon on the right.


Departamento de Patologia e Oncologia
Departamento https://sigarra.up.pt/fmup/pt/uni_geral.unidade_view?pv_unid...
Faculdade de Medicina da Universidade do Porto (FMUP)
Responsável: Manuel Alberto Coimbra Sobrinho Simões


Departamento de Física e Astronomia
Departamento https://sigarra.up.pt/fcup/pt/uni_geral.unidade_view?pv_unid...
Faculdade de Ciências da Universidade do Porto (FCUP)
Responsável: Orfeu Bertolami Neto

Departamento de Geografia
Departamento https://sigarra.up.pt/flup/pt/uni_geral.unidade_view?pv_unida...
Faculdade de Letras da Universidade do Porto (FLUP)
Responsável: Teresa Maria Vieira de Sá Marques


Departamento de Apoio ao Estudante
Departamento https://sigarra.up.pt/sasup/pt/uni_geral.unidade_view?pv_uni...
Serviços de Ação Social da Universidade do Porto (SASUP)
Responsável: Maria Cristina Sampaio da Mota e Silva

Type query. The intention of the query is to find entities of a given type or class.






[Todos](#)
[Salas](#)
[Pessoal](#)
[Estudantes](#)
[Notícias](#)
[Ferramentas de Pesquisa](#)

4 resultados (3.81 segundos)


José Luís da Silva Devezas

Sala
I123




José Luís da Silva Devezas

Estudante
https://sigarra.up.pt/flup/pt/vld_entidades_geral.entidade_pa...


Faculdade de Belas Artes da Universidade do Porto (FBAUP) (mais 4)

Curso: Doutoramento em Informática

Código: 200303288




José Luís da Silva Devezas


Pessoal

https://sigarra.up.pt/feup/pt/func_geral.FormView?P_CODIGO...

Faculdade de Engenharia da Universidade do Porto (FEUP) (mais 1)

Código: 493720 **Sigla:** JLSD

 1317

I123

Sala

https://sigarra.up.pt/feup/pt/instal_geral.espaco_view?pv_id=...

Laboratório de I&D de Sistemas de Informação, Faculdade de Engenharia da Universidade do Porto (FEUP)

Responsáveis: [Carla Alexandra Teixeira Lopes](#)

Edifício: Electrotecnia (I) **Andar:** 1

Attribute query. The intention of the query is to find values for a given attribute of a particular entity or type.



6 resultados obtidos em 2,88 segundos

José Luís da Silva Devezas

Faculdade de Belas Artes da Universidade do Porto (FBAUP), Faculdade de Letras da Universidade do Porto (FLUP), Faculdade de Engenharia da Universidade do Porto (FEUP), Faculdade de Economia da Universidade do Porto (FEP), Faculdade de Ciências da Universidade do Porto (FCUP)

Estudante

score = 0.002238

**Tiago Nuno Mesquita Folgado leitão Devezas**

Faculdade de Engenharia da Universidade do Porto (FEUP) Domingo, 07 Outubro 2018, 00h32

Pessoal

score = 0.001971

**José Luís da Silva Devezas**

Faculdade de Engenharia da Universidade do Porto (FEUP), Faculdade de Ciências da Universidade do Porto (FCUP) Domingo, 02 Setembro 2018, 00h26

Pessoal

score = 0.001106

**Tiago Nuno Mesquita Folgado Leitão Devezas**

Faculdade de Belas Artes da Universidade do Porto (FBAUP), Faculdade de Letras da Universidade do Porto (FLUP), Faculdade de Engenharia da Universidade do Porto (FEUP), Faculdade de Economia da Universidade do Porto (FEP), Faculdade de Ciências da Universidade do Porto (FCUP)

Estudante

score = 0.001789

**I123**

Faculdade de Engenharia da Universidade do Porto (FEUP)

Sala

score = 0.000431

**Projetos no laboratório SAPO/U.Porto**

Faculdade de Engenharia da Universidade do Porto (FEUP) Sexta-feira, 28 Junho 2013, 00h00

Notícia

score = 0.0



<<

Anterior

Página 1 de 1

Seguinte

>>

Entidades

Tiago Nuno Mesquita Folgado leitão Devezas



José Luís da Silva Devezas

Ligações

Faculdade: Faculdade de Engenharia, Faculdade de Engenharia da Universidade do Porto

Sala: I123

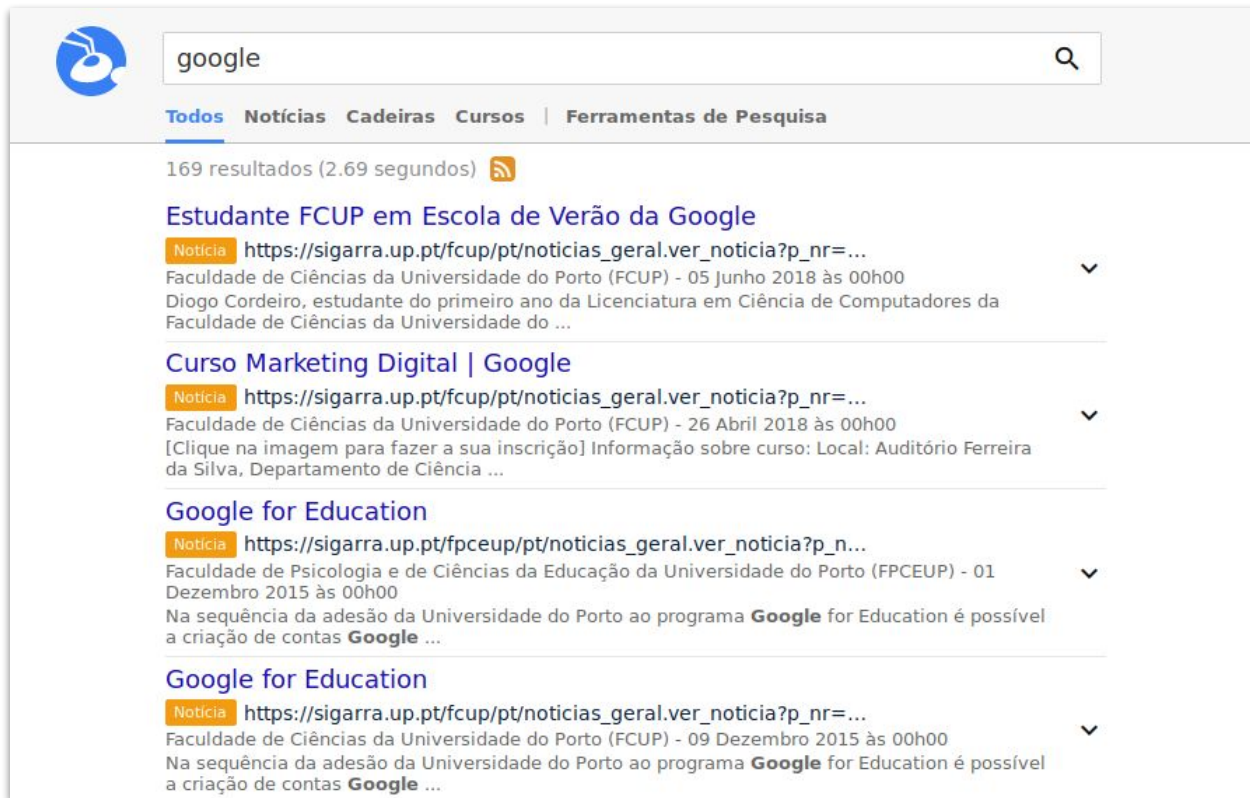
Projetor Não

Edifício Electrotecnia (I)

Mapa https://sigarra.up.pt/feup/pt/instal_neral?nei_mana?nv_id=77486

[Ver mais](#)

Relation query. The intention of the query is to find how two or more entities or types are related.



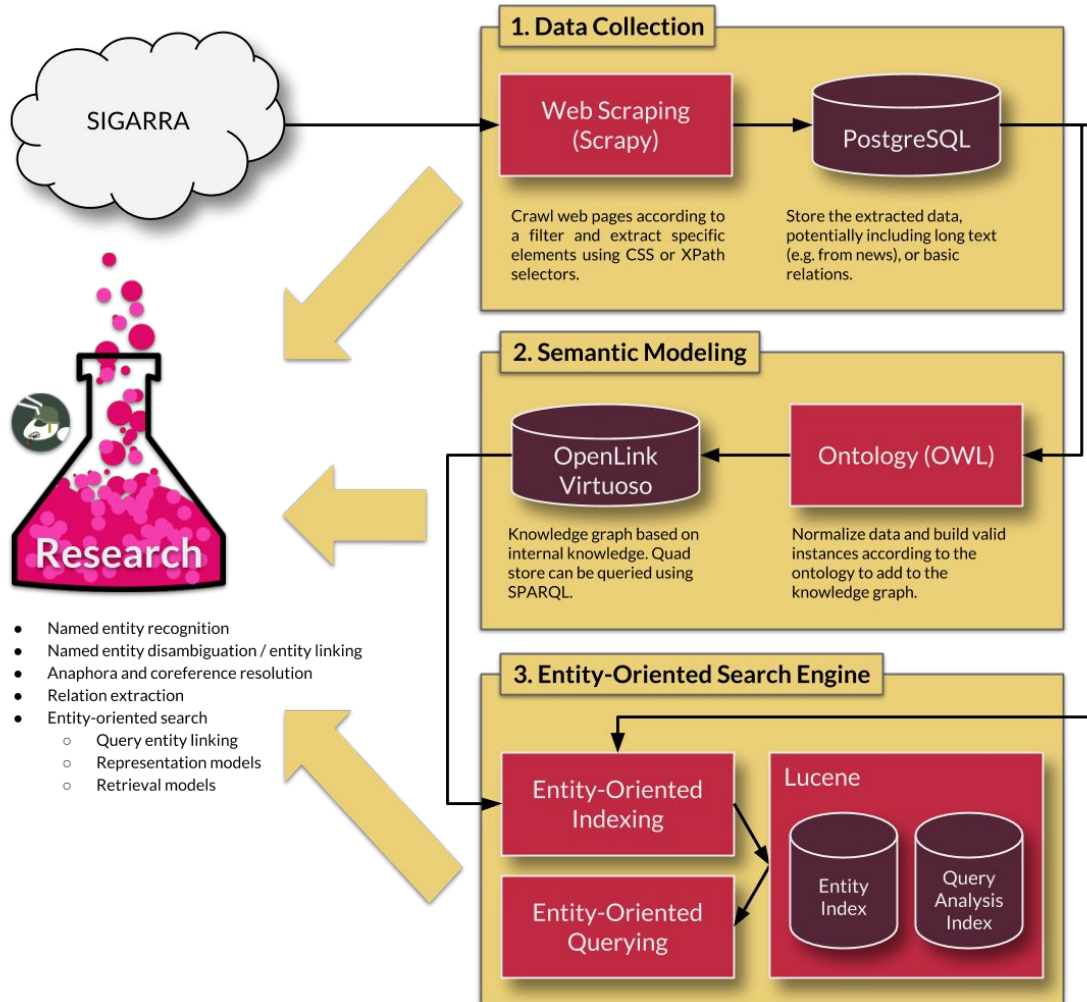
The screenshot shows a Google search interface with the search bar containing the word "google". Below the search bar, there are navigation links: "Todos", "Notícias", "Cadeiras", "Cursos", and "Ferramentas de Pesquisa". The search results show 169 results in 2.69 seconds. The first four results are news items from the University of Porto (FCUP) website, all with the URL "https://sigarra.up.pt/fcup/pt/noticias_geral.ver_noticia?p_nr=...". Each result has a "Notícia" label and a dropdown arrow. The results are:

- Estudante FCUP em Escola de Verão da Google**
Faculdade de Ciências da Universidade do Porto (FCUP) - 05 Junho 2018 às 00h00
Diogo Cordeiro, estudante do primeiro ano da Licenciatura em Ciência de Computadores da Faculdade de Ciências da Universidade do ...
- Curso Marketing Digital | Google**
Faculdade de Ciências da Universidade do Porto (FCUP) - 26 Abril 2018 às 00h00
[Clique na imagem para fazer a sua inscrição] Informação sobre curso: Local: Auditório Ferreira da Silva, Departamento de Ciência ...
- Google for Education**
Faculdade de Psicologia e de Ciências da Educação da Universidade do Porto (FPCEUP) - 01 Dezembro 2015 às 00h00
Na sequência da adesão da Universidade do Porto ao programa **Google** for Education é possível a criação de contas **Google** ...
- Google for Education**
Faculdade de Ciências da Universidade do Porto (FCUP) - 09 Dezembro 2015 às 00h00
Na sequência da adesão da Universidade do Porto ao programa **Google** for Education é possível a criação de contas **Google** ...

Keyword query. Anything that doesn't fit the previous four categories.

Search engine architecture

ANT components, from data collection to search.



1. Data Collection

Scrapy

I 123

Edifício: Electrotecnia (I)

Utilização: Laboratório - Investigação

Área (m²): 60

Telefone: 1317

Responsáveis: • Carla Alexandra Teixeira Lopes

Ocupantes: • Inês Dias Koch
• Joana Patrícia de Sousa Rodrigues
• João Daniel Aguiar de Castro
• João Miguel Rocha da Silva
• José Luís da Silva Devesas
• Nelson Miguel da Costa Martins Pereira
• Tiago Nunes Mesquita Folgado Leitão Devesas



Spiders

Pipelines

Models

Tables:
One per entity +
auxiliary tables

Materialized view:
lucene_documents

Specify:

- Allowed domains (up.pt)
- List of start URLs (all SIGARRAs)
- LinkExtractor rules

Crawl:

- Launch crawlers from start URLs
- Fetch and parse web pages
- Extract URLs from anchors
- URLs matching rules will be scraped
- Follow uncrawled up.pt URLs

Scrape:

- Manually define XPath and CSS selectors for elements to be stored in the model.

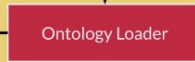
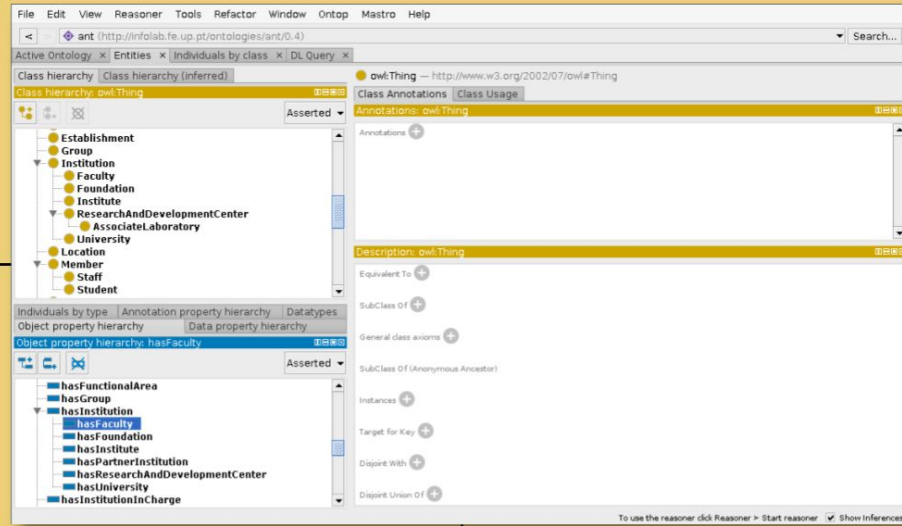
Transform an HTTP response into an **item** through predefined CSS and XPath selectors.

Transform an **item** from a spider into an **object** from the a model, and deal with storage operations.

e.g., object-relational mapping (ORM)

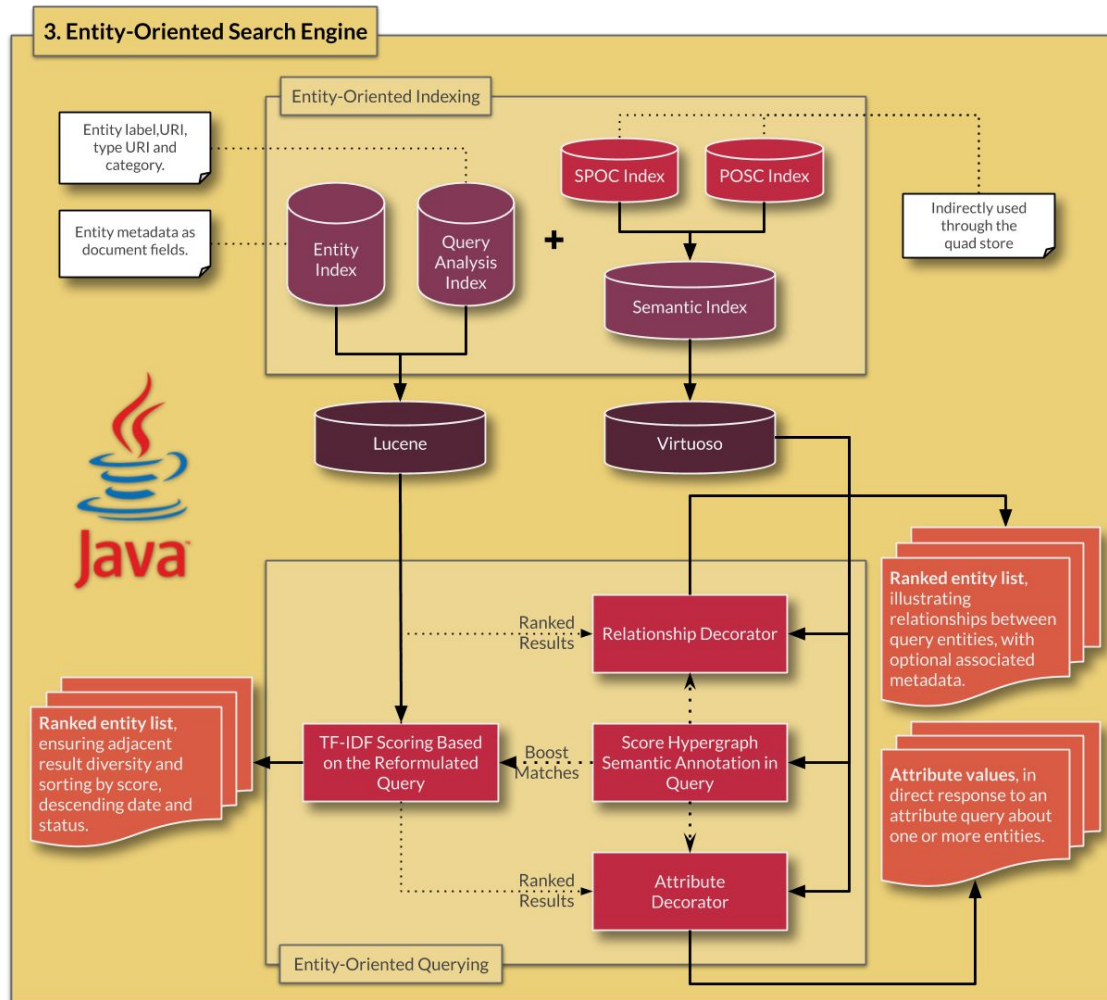
PostgreSQL

2. Semantic Modeling



Transforms PostgreSQL relational data into RDF triples to store in Virtuoso.

3. Entity-Oriented Search Engine





REST API

- ANT provides access to search-related services via a REST API.
- We use the OpenAPI 2.0 format (Swagger) to document the API.
 - <http://ant.fe.up.pt/api/>
 - <https://swagger.io/specification/>
- Which makes it possible to easily provide a console for API exploration.
 - <http://ant.fe.up.pt/api-console/>

- Supported services are classified into six categories:
 - Analytics
 - Autocomplete
 - Decorators*
 - JavaScript
 - Log
 - Search*

* Critical services.

Conclusions

Final remarks and related projects.





Final remarks

- The ANT search engine is serving the local academic community, despite infrastructure and human resource limitations (it's a prototype).
- At the same time, it collects implicit relevance feedback, based on result clicks for issued queries.
- ANT is also a platform of collaboration for multiple areas of research:
 - Web Design
Collaboration with MM for the development of the front-end.
 - User Experience
MM dissertation in entity-oriented search interfaces.
 - Information Extraction
MIEIC dissertation in named entity recognition for portuguese web text.



Related projects

Army ANT



- Serving the research needs in the area of entity-oriented search.
- Supporting the study of innovative ideas in search, providing tools for exploration and evaluation.

PhD thesis



- “Graph-Based Entity-Oriented Search”
 - Joint representation of text, entities and their relations.
 - Generalization of entity-oriented search tasks.
 - Improvement of search effectiveness?
- Exploration of random walks in graphs and hypergraphs.

Thank you!

<https://ant.fe.up.pt>

José Devezas is supported by research grant PD/BD/128160/2016, provided by the Portuguese national funding agency for science, research and technology, Fundação para a Ciência e a Tecnologia (FCT), within the scope of Operational Program Human Capital (POCH), supported by the European Social Fund and by national funds from MCTES.
