# ANT – R&D Platform for Entity-Oriented Search

IEEE SYP 2018 Porto

**Last update:** Jul 26, 2018

José Devezas    MAP-i 2016/2017    INESC TEC & FEUP InfoLab

**Supervisor:**    Sérgio Nunes    U.Porto
**External advisor:**    Bruno Martins    U.Lisboa
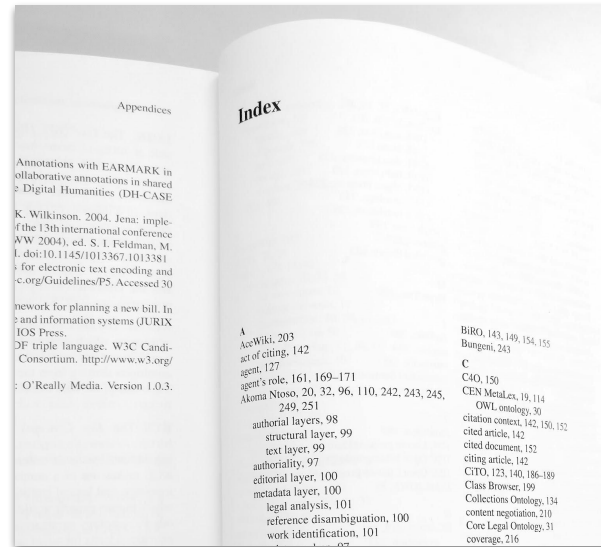
Universidade do Minho

universidade de aveiro

U.PORTO

# Contents

# Introduction

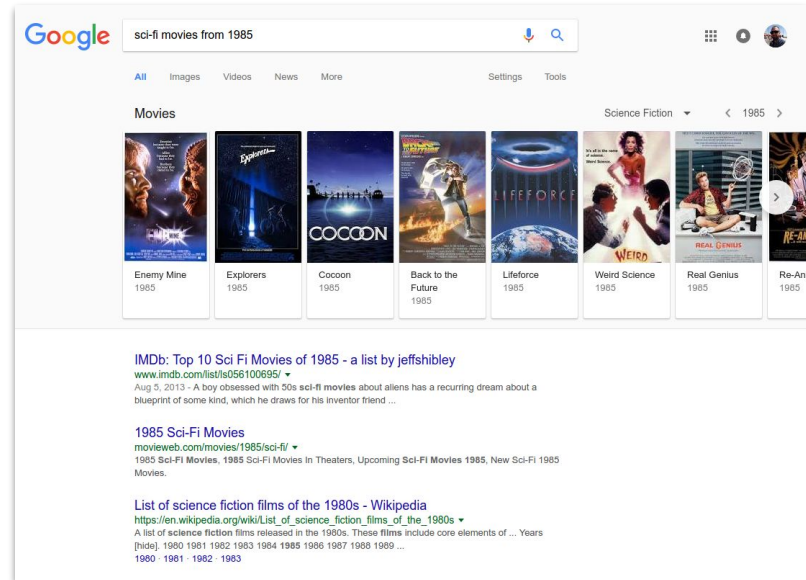**What is entity-oriented search and why does it matter?**

# **Keyword-based** vs entity-oriented search

- Keyword-based search was modeled after the back-of-the-book index.
- Finding relevant content involved:
    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.

# Keyword-based vs <u>entity-oriented</u> search

- Entity-oriented search makes use of:
  - Natural language understanding:
    - For queries;
    - And documents.
  - Structured data from knowledge bases.
- Making it possible to answer queries like:
  - [ sci-fi movies from 1985 ]
- By returning a combination of:
  - Text documents;
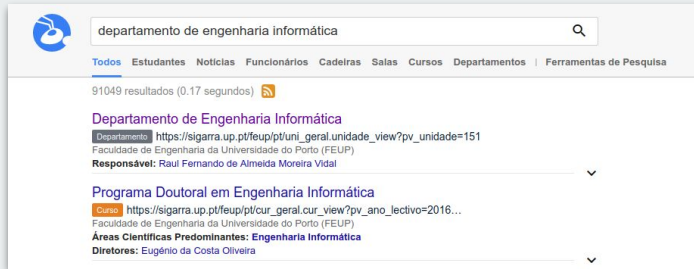  - And entities (e.g., movies).



5

# The relevance of entities in search

- In queries:.
  - A study of the AOL Query Log showed that:
    - 18-39% queries directly reference entities;
    - 73-87% queries contain at least one entity.

- In documents:
  - The annotated CoNLL 2003 English training set contained:
    - 14,987 sentences;
    - 23,499 entities;
    - Resulting in 1.6 entities per sentence.

# ANT

**Searching for information at the University of Porto.**

ANT



departamento de engenharia informática

Todos  Estudantes  Noticias  Funcionários  Cadeiras  Salas  Cursos  Departamentos | Ferramentas de Pesquisa

91049 resultados (0.17 segundos)

Departamento de Engenharia Informática
Departamento https://sigarra.up.pt/feup/pt/uni_geral.unidade_view?pv_unidade=151
Faculdade de Engenharia da Universidade do Porto (FEUP)
Responsável: Raul Fernando de Almeida Moreira Vidal

Programa Doutoral em Engenharia Informática
Curso https://sigarra.up.pt/feup/pt/cur_geral.cur_view?pv_ano_lectivo=2016…
Faculdade de Engenharia da Universidade do Porto (FEUP)
Áreas Científicas Predominantes: Engenharia Informática
Diretores: Eugénio da Costa Oliveira

Ad hoc search of eNtities and Text.

- ANT is an entity-oriented search engine, built to support the five query categories defined by Pound et al. (2010):
  - Entity query;
  - Type query;
  - Attribute query;
  - Relation query;
  - Keyword query.
- It is supported by two Lucene indexes:
  - Query analysis index;
  - Entity index.
- And a Virtuoso RDF triplestore:
  - Useful for relation queries.

# How does ANT understand queries?

- Query segmentation based on the retrieval of matching entities for all query $n$-grams up to a maximum value of $n$.
- Semantic tagging of query segments based on the probability of associating a given type of entity to an $n$-gram.

**Query:** josé sérgio sobral nunes informática

| $n = 1$ | $n = 2$ | $n = 3$ | $n = 4$ |
|---|---|---|---|
| josé | josé sérgio | josé sérgio sobral | josé sérgio sobral nunes |
| sérgio | sérgio sobral | sérgio sobral nunes | sérgio sobral nunes informática |
| sobral | sobral nunes | sobral nunes informática | |
| nunes | nunes informática | | |
| informática | | | |

Devezas, J. and S. Nunes (2016). Index-Based Semantic Tagging for Efficient Query Interpretation. In Proceedings of the 6th International Conference of the CLEF Initiative (CLEF 2016), Évora, Portugal.

# How does ANT understand queries?

- The **actual method** we ended up using is a variation of this that we called "**Score Hypergraph**".
  - TF-IDF scores instead of probabilities.
  - Dedicated query analysis index to search for entities matching $n$-grams.
  - Hypergraph* of $n$-grams to resolve query segment overlaps and to fix bugs with the previous approach.

\* A hypergraph is a generalization of a graph, where edges can have an arbitrary number of nodes.

     Think of a social network modeling binary friendship relations (a social graph).

     And then think of a social network that also models groups of multiple friends (a social hypergraph).

10

# Score Hypergraph

**Query segmentation and semantic tagging in ANT.**

Note: Displayed weights are the TF-IDF of the top-1 matched entity.
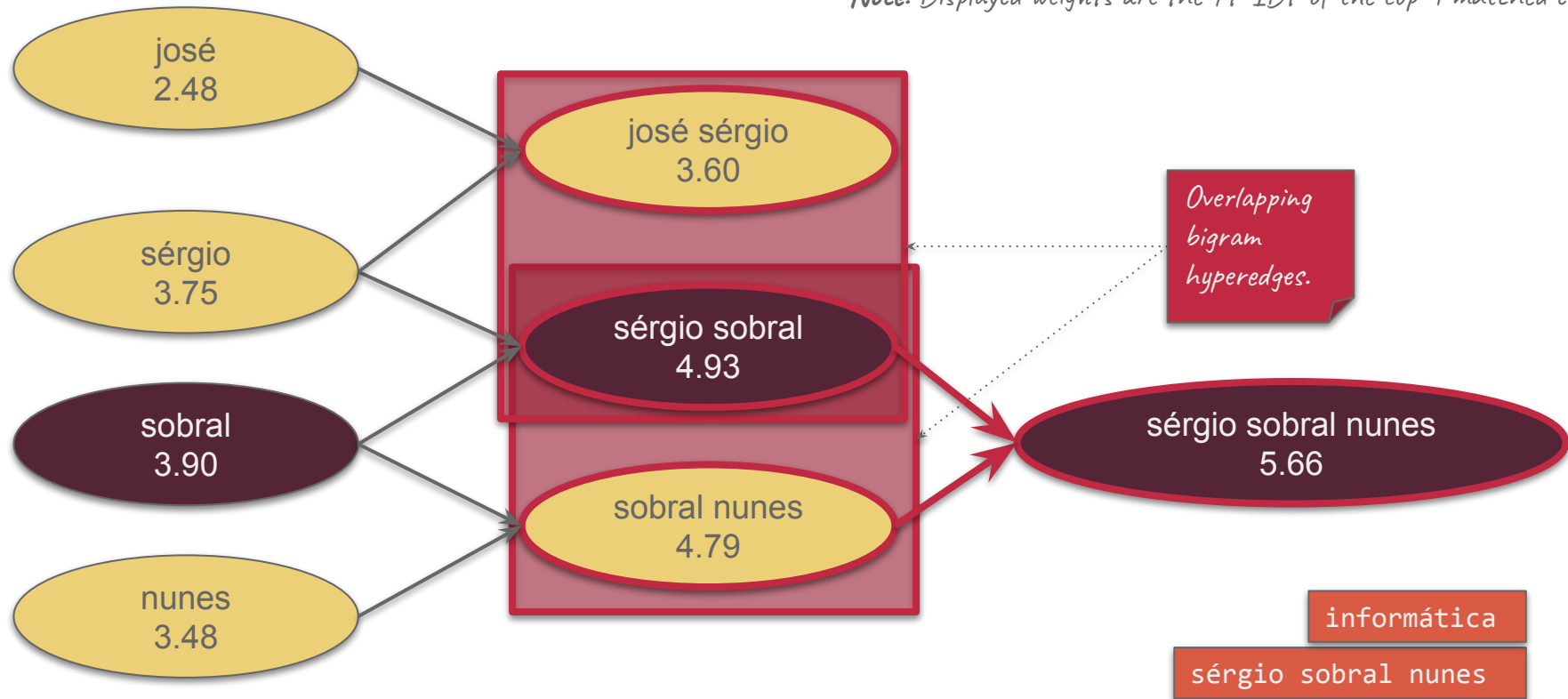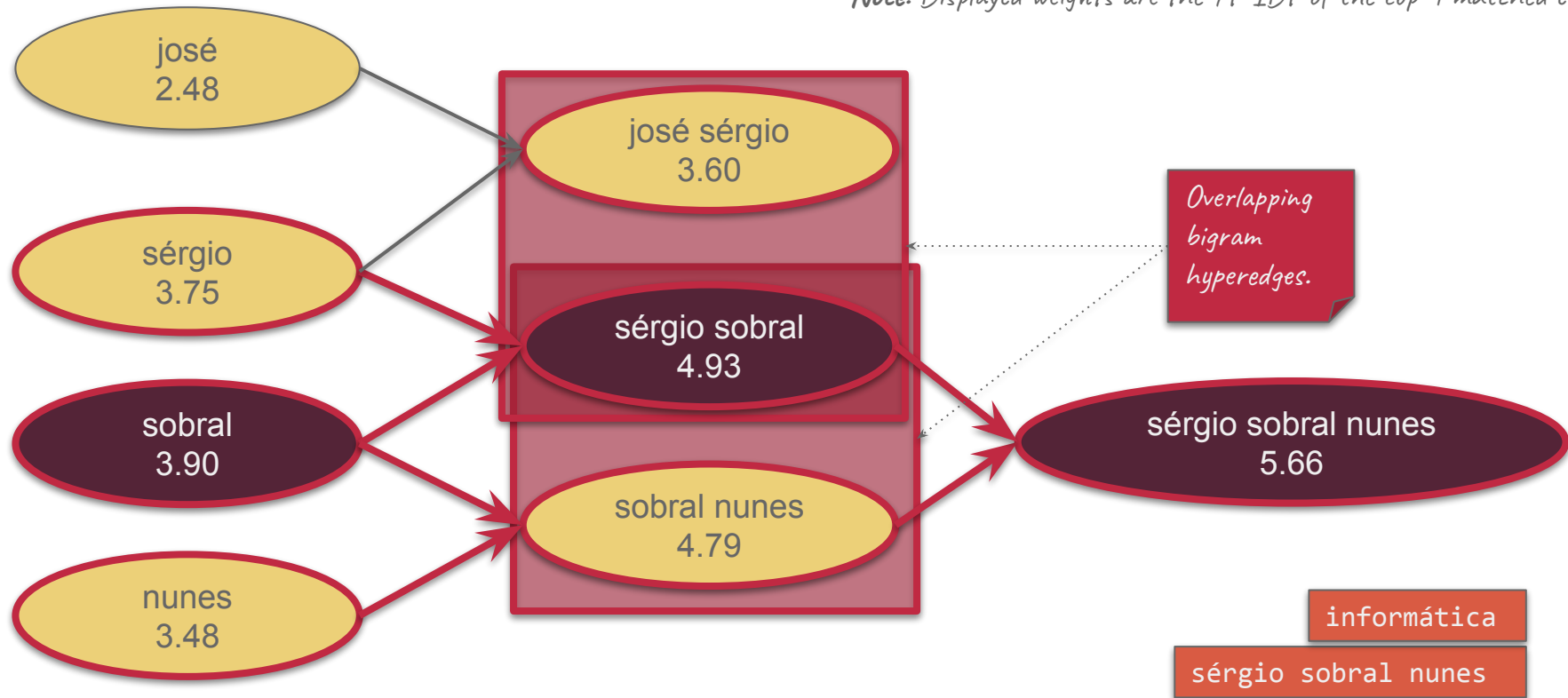
Overlapping bigram hyperedges.

Query: josé sérgio sobral nunes informática

**Note:** Displayed weights are the TF-IDF of the top-1 matched entity.

josé
2.48

sérgio
3.75

sobral
3.90

nunes
3.48

josé sérgio
3.60

sérgio sobral
4.93

sobral nunes
4.79

sérgio sobral nunes
5.66

Overlapping bigram hyperedges.

informática

sérgio sobral nunes

**Query:** josé sérgio sobral nunes informática

**Note:** Displayed weights are the TF-IDF of the top-1 matched entity.

josé
2.48

sérgio
3.75

sobral
3.90

nunes
3.48

josé sérgio
3.60

sérgio sobral
4.93

sobral nunes
4.79

Overlapping bigram hyperedges.

sérgio sobral nunes
5.66

informática

sérgio sobral nunes

**Query:** josé sérgio sobral nunes informática

**Note:** Displayed weights are the TF-IDF of the top-1 matched entity.

josé
2.48

sérgio
3.75

sobral
3.90

nunes
3.48

josé sérgio
3.60

sérgio sobral
4.93

sobral nunes
4.79

Overlapping bigram hyperedges.

sérgio sobral nunes
5.66

informática

sérgio sobral nunes

**Query:** josé sérgio sobral nunes informática

josé
2.48

informática

sérgio sobral nunes

**Query:** josé sérgio sobral nunes informática

josé
2.48

informática

sérgio sobral nunes

**Query:** josé sérgio sobral nunes informática

josé
2.48

informática

sérgio sobral nunes

josé

**Query:** josé sérgio sobral nunes informática

informática

sérgio sobral nunes

josé

**Query:** josé sérgio sobral nunes informática

**Query:** josé sérgio sobral nunes informática

1. The query was segmented based on the $n$-grams with the highest-scoring entities.

**Query:** josé sérgio sobral nunes informática

Staff      Staff      Department

2. The query was assigned semantic tags based on the type of the highest-scoring entity.

3. From the semantic tag, we directly derived a higher level tag that could either be ENTITY (e.g., instance of Staff class), ATTRIBUTE (e.g., property) or TYPE (e.g., Staff class).

4. Based on the combination of higher level tags, we conditionally obtained the query category.

# Army ANT

**Researching entity-oriented search.**

# Army ANT

Army ANT The altruistic ant

Search   Evaluation   About

musician                                    **Search**    Learn mode

INEX 3T-NL - Hypergraph-o ⬍   Random Walk Score ⬍   l   3 ⬍   r   10 ⬍

3 results (0.25 seconds)

Henry Kaiser (musician)

Warren Smith (jazz musician)

A workbench for innovation in entity-oriented search.

- Indexing unit: documents with $doc\_id$, $text$ and $triples$.
- Able to define **readers** that work as iterators of documents.
- Able to implement retrieval models **(engines)** by implementing $index()$ and $search()$ methods.

# Army ANT

A workbench for innovation in entity-oriented search.

Front-end provides:

- Standard **search interface**, where you can select an index and a ranking function.
- **Learn mode** interface with:
  - Results without metadata;
  - Score component visualization;
  - Trace for the active query;
  - Ranking function details;
  - Collection description.
- **Evaluation interface** supporting:
  - Topics+Assessments (INEX Ad Hoc and INEX XML Entity Ranking);
  - Topics (TREC Common Core);
  - Living Labs API (TREC OpenSearch).

# Command Line Interface

`index` | `search` | `inspect` | `analysis` | `sampling` | `features` | `extras` | `evaluation` | `server`

# ./army-ant.py

➜ `index`
  ◆ Index a supported collection (i.e., based on an implemented reader), using one of the available engines.
➜ `search`
  ◆ Search one of the supported indexes (has an interactive mode to avoid preload latency).
➜ `inspect`
  ◆ Extract several features from a particular index (supported features depend on the engine).

➜ `analysis`
  ◆ `rws-rank-concordance`
    ● Analyze rank concordance for Random Walk Score (Hypergraph-of-Entity).
➜ `sampling`
  ◆ Create a subset of one of the supported collections.
➜ `features`
  ◆ Extract features (usually from a collection), such as word embeddings and similarities.

# ./army-ant.py

➔ `extras`
  ◆ `fetch-wikipedia-images`
    ● Obtain the Wikipedia image URL for documents stored in the database.
  ◆ `word2vec-knn`
    ● Return a ranked list of the $k$-nearest neighbors for a given word.
  ◆ `word2vec-sim`
    ● Measure the similarity between the embeddings for two words.

➔ `evaluation`
  ◆ Queue and run an evaluation task using a supported evaluator (for now, it supports INEX and Living Labs API; soon it will support TREC qrels).

➔ `server`
  ◆ Launch the web server with a search interface, a learn mode and an evaluation panel.

35

# Configuration

**Based on YAML and mostly used by the web interface.**

# Example file

Global settings for metadata storage, evaluation metrics and location, and reserved heap space for Java-based engines.

```
defaults:
  db:
    location: mongo
    name: army_ant
    type: mongo
  eval:
    metrics:
      favorite:
        - GMAP
        - MAP
        - NDCG@10
        - P@10
    location: /home/army-ant/data/eval
  service:
    ner:
      entity_list: /home/army-ant/data/people.txt
  depend:
    stanford-ner: /opt/stanford-ner-2015-12-09
  jvm:
    memory: 5120
    other_args: -XX:+UseConcMarkSweepGC
```

# Example file

Lucene index and ranking functions configuration.

```
engines:
  lucene-inex-3t-nl:
    name: INEX 3T-NL - Lucene
    db:
      name: inex
    index:
      type: lucene
      location: /home/army-ant/data/indexes/lucene
    ranking:
      default:
        id: tf_idf
      functions:
        tf_idf:
          name: TF-IDF
        bm25:
          name: BM25
          params:
            k1: [1.2, 1, 1.8]
            b: [0.75, 0.5, 1]
        dfr:
          name: DFR
          params:
            BM: [BE, G, P, D, In, Ine, IF]
            AE: [L, B, Disabled]
            N: [H1, H2, H3, Z, Disabled]
```

# Example file

Hypergraph-of-entity index and ranking functions configuration.

```
engines:
  hgoe-inex-3t-nl:
    name: INEX 3T-NL - Hypergraph-of-Entity
    db:
      name: inex
    index:
      type: hgoe
      location: /home/army-ant/data/indexes/hgoe
      preload: true
    ranking:
      default:
        id: random_walk
        params:
          l: 2
          r: 10
      functions:
        jaccard:
          name: Jaccard Score
        random_walk:
          name: Random Walk Score
          params:
            l: [1, 2, 3, 4, 5, 6]
            r: [10, 25, 50, 100, 1000]
```

# Web Interface

**Searching, learning and evaluating.**

1. **Search interface**, showing a query over an hypergraph-of-entity index, using the random walk score as the ranking function, with $\ell$ = 3 and $r$ = 10.

There's something new over here, that we will show in the demo.

1. **Search interface**, showing a query over an hypergraph-of-entity index, using the random walk score as the ranking function, with $\ell$ = 3 and $r$ = 10.

**Learn Mode**
EXPLORE, UNDERSTAND, ANALYZE

| Results | Trace | Model | Collection |
|---------|-------|-------|------------|

| Rank | Score(q, d) | Doc ID |
|------|-------------|--------|
| 1 | 0.244444 | 9934261 |
| 2 | 0.111111 | 1193582 |
| 3 | 0.083333 | 16265226 |
| 4 | 0.033333 | 19127472 |
| 5 | 0.011111 | 3764544 |

2. **Learn mode**, showing the results with the ranking, score and document ID.

3. **Learn mode**, showing a trace ("instanced explain") and respective ASCII export for the hypergraph-of-entity engine.

4. **Learn mode**, showing the score components visualization based on the parallel coordinates system. Displayed score components are based on the graph-of-entity and the entity weight ranking function.

# 5. Learn mode: model

Illustrated with the description of the hypergraph-of-entity, dynamically showing which index extensions are enabled and which parameter values were used in the active query.

# 6. **Learn mode:** collection

Shows the description of the indexed collection, including the source of the data, the temporal coverage, a free text description, an example of the data and a reference paper.



**Learn Mode**
EXPLORE, UNDERSTAND, ANALYZE

| Results | Trace | Model | Collection |

**Source**
INEX Document Collections - Adhoc Track (2009-2010) and Wikipedia Collection (2009-)

**Date**
October 8, 2008

**Description**
"Starting in 2009, INEX uses a new document collection based on the Wikipedia. The original Wiki syntax has been converted into XML, using both general tags of the layout structure (like article, section, paragraph, title, list and item), typographical tags (like bold, emphatic), and frequently occurring link-tags. The annotation is enhanced with semantic markup of articles and outgoing links, based on the semantic knowledge base YAGO, explicitly labeling more than 5,800 classes of entities like persons, movies, cities, and many more. For a more technical description of a preliminary version of this collection, see [9].

The collection was created from the October 8, 2008 dump of the English Wikipedia articles and incorporates semantic annotations from the 2008-w40- 2 version of YAGO. It contains 2,666,190 Wikipedia articles and has a total uncompressed size of 50.7 Gb. There are 101,917,424 XML elements of at least 50 characters (excluding white-space)."

**Example**

```
<?xml version="1.0" encoding="UTF-8"?>
<!-- generated by CLiX/Wiki2XML [MPI-Inf, MMCI@UdS] $LastChangedRevision: 92 $ on 16.04.2009
```

```
<b>Johnny Burke</b> (1851 - 1930) was a <country wordnetid="108544813" confidence="0.9508927676800064">
<link xlink:type="simple" xlink:href="../561/697561.xml">
Newfoundland</link></country>
songwriter and musician. He was nicknamed the 'Bard of Prescott Street'. He wrote many popular songs that
artists in the 1930s and 1940s released.</p>
```

**Paper**
Geva S., Kamps J., Lethonen M., Schenkel R., Thom J.A., Trotman A. (2010) Overview of the INEX 2009 Ad Hoc Track. In: Geva S., Kamps J., Trotman A. (eds) Focused Retrieval and Evaluation. INEX 2009. Lecture Notes in Computer Science, vol 6203. Springer, Berlin, Heidelberg

*Springer Link*

47

7. **Evaluation module**, showing the task launching form.

8. **Evaluation module**, showing a finished task and its expanded results.

9. **Evaluation module**, showing the global evaluation export, for comparing the results of multiple parameter configurations among different tasks.

# Conclusions

**Final remarks and interactive demo.**

# Final remarks

- The ANT search engine is serving the local academic community and giving us a test platform.

- Army ANT is serving the research needs in the area and supporting my PhD.

- By the way, my thesis topic is "Graph-Based Entity-Oriented Search":
  - If you're interested on the topic, feel free to look me up on ANT and contact me.
  - I'm also exploring hypergraphs as an alternative, higher-level, representation model.
  - The goal is to integrate text and knowledge in a joint model.
  - And to provide a generalized model to support entity-oriented search tasks.

# Installing Army ANT demo using Docker

- First install Docker Compose:
    - https://docs.docker.com/compose/install/

- And then get Army ANT install repository:
    - https://github.com/feup-infolab/army-ant-install

- Follow the instructions  to launch Army ANT and explore the included Lucene and hypergraph-of-entity indexes:
    - git clone git@github.com:feup-infolab/army-ant-install.git
    - cd army-ant-install
    - git checkout ieee-syp-2018
    - docker-compose up

# Thank you!

**https://ant.fe.up.pt**
**https://github.com/feup-infolab/army-ant**