

ANT Searching for Information at the University of Porto

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
 - The relevance of entities in search
 - Entity-oriented search
- ANT
 - Search engine architecture
 - Query understanding
 - Score Hypergraph
- Conclusions
 - Final remarks
 - Related projects
 - A vision for the future

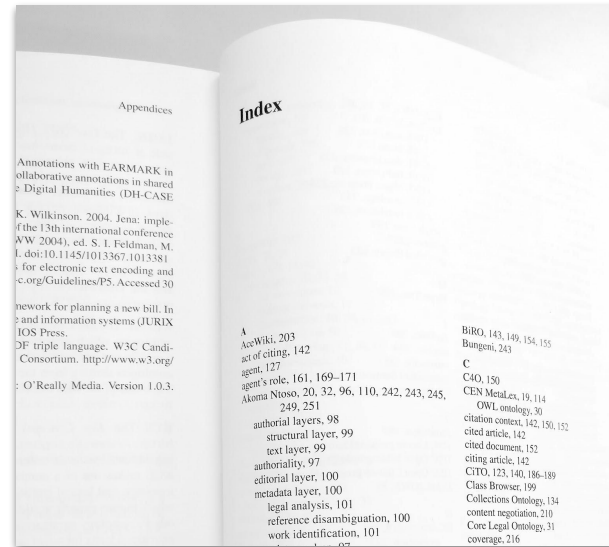
Introduction

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



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 the user's information need.
- Requires interpretation of query meaning and document semantics.
- Combines unstructured data from text and structured data from knowledge bases.
 - Google Knowledge Graph
 - Google Knowledge Vault
 - DBpedia
 - Wikidata

The screenshot shows a Google search for "president of portugal 2018". The search results page displays the Knowledge Graph panel for Marcelo Rebelo de Sousa, the current President of Portugal. The panel includes a portrait of Marcelo Rebelo de Sousa, his name, and a list of "People also search for" with small portraits and names: António Costa, Eduardo Ferro Rodrigues, Aníbal Cavaco Silva, José Sócrates, Pedro Passos Coelho, Rui Rio, and Jorge Sampaio. Below the panel, there is a link to the Wikipedia page for "President of Portugal" and a snippet of text stating: "The current President of Portugal is **Marcelo Rebelo de Sousa**, who took office on 9 March 2016." Additional information includes: "Term length: Five years; Renewable once, con...", "Residence: Belém Palace", "Salary: €93,364.74 (2015); (€6,668.91/month)", and "First holder: Manuel de Arriaga".

Search using semantic matching

- It becomes possible to, more adequately, answer queries like:
 - [president of portugal 2018]
 - [sci-fi movies from 1985]
- Instead of documents containing these keywords;
- An entity or list of entities is directly provided as the answer.
- Avoids the hassle of having to skim through documents to find the answer.

Google sci-fi movies from 1985

All Images Videos News More Settings Tools

Movies Science Fiction < 1985 >

Enemy Mine 1985 Explorers 1985 Cocoon 1985 Back to the Future 1985 Lifeorce 1985 Weird Science 1985 Real Genius 1985 Re-Ani 1985

IMDb: Top 10 Sci Fi Movies of 1985 - a list by jeffshibley
www.imdb.com/list/ls056100695/ ▼
 Aug 5, 2013 - A boy obsessed with 50s sci-fi movies about aliens has a recurring dream about a blueprint of some kind, which he draws for his inventor friend ...

1985 Sci-Fi Movies
movieweb.com/movies/1985/sci-fi/ ▼
 1985 Sci-Fi Movies, 1985 Sci-Fi Movies In Theaters, Upcoming Sci-Fi Movies 1985, New Sci-Fi 1985 Movies.

List of science fiction films of the 1980s - Wikipedia
https://en.wikipedia.org/wiki/List_of_science_fiction_films_of_the_1980s ▼
 A list of science fiction films released in the 1980s. These films include core elements of ... Years [hide] 1980 1981 1982 1983 1984 1985 1986 1987 1988 1989 ... 1980 1981 1982 1983



The relevance of entities in search

In queries:

- A 2007 study of the AOL Query Log showed that:
 - 18-39% queries directly refer to entities;
 - 73-87% queries contain at least one entity.

In documents:

- The annotated CoNLL 2003 English training set contains:
 - 14,987 sentences;
 - 23,499 entities;
 - 1.6 entities per sentence.



Entity-oriented search

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.

In ANT, entities can be:

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

ANT

Search engine architecture and query understanding method.



The logo for ANT, consisting of the letters 'ANT' in a bold, black, sans-serif font. Above the letters is a horizontal bar with a green segment on the left and an orange segment on the right.

Ad hoc search of eNtities and Text.

- Supports the five query categories 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.

Logo: 

Search bar: 

Navigation: [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.*

departamento

Todos Departamentos | Ferramentas de Pesquisa

82 resultados (2.12 segundos)


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

Noticia

score = 0.0



Anterior

Página 1 de 1

Seguinte



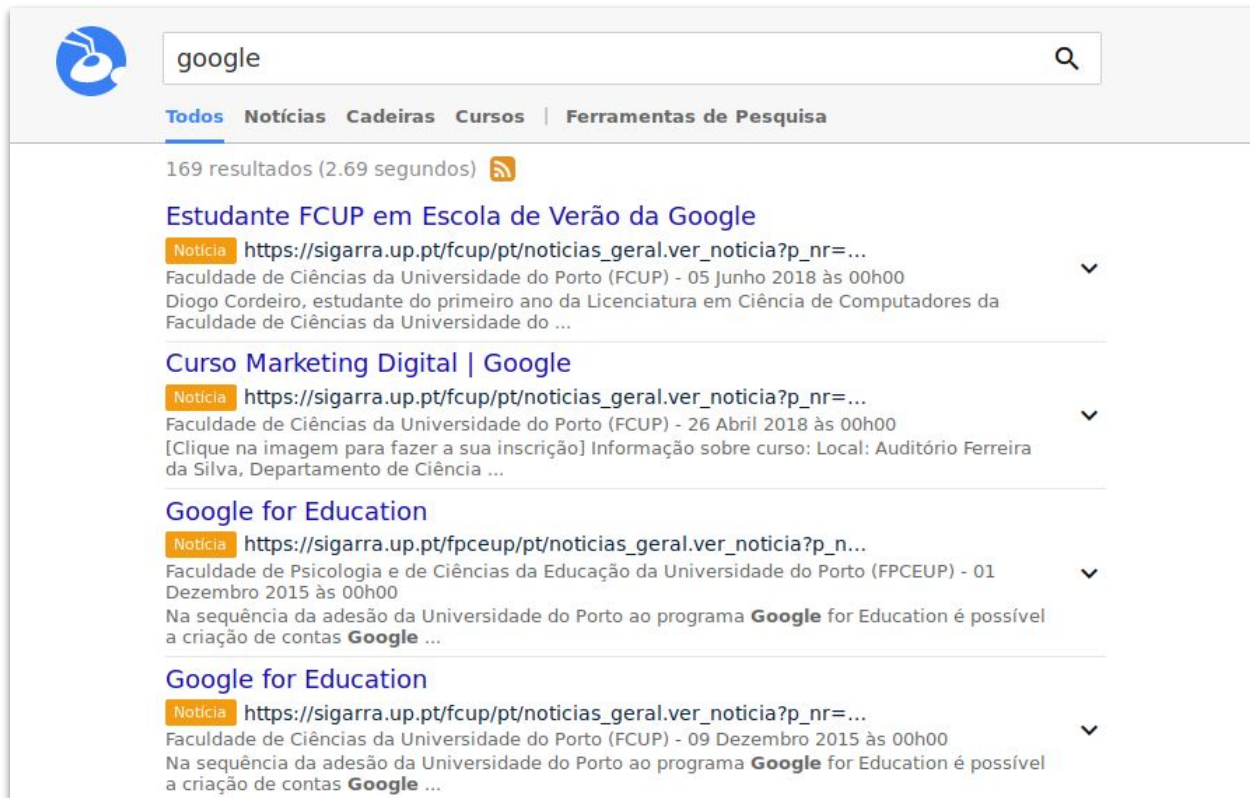
Entidades

**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.



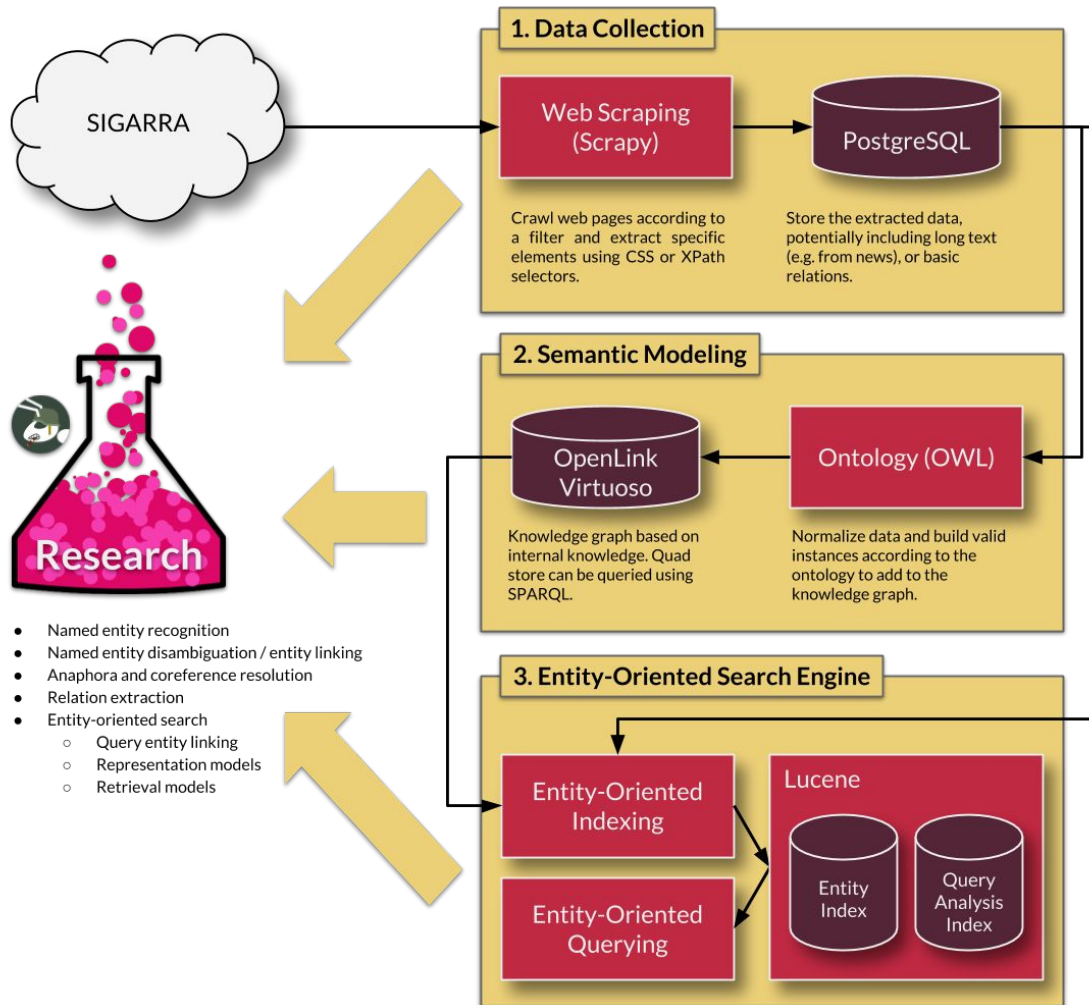
The image shows a search engine interface with the following elements:

- Search Bar:** Contains the text "google" and a magnifying glass icon.
- Navigation:** A horizontal menu with tabs: "Todos", "Notícias", "Cadeiras", "Cursos", and "Ferramentas de Pesquisa".
- Results Summary:** "169 resultados (2.69 segundos)" followed by a feed icon.
- Result 1:**
 - Title:** "Estudante FCUP em Escola de Verão da Google"
 - Snippet:** "Notícia https://sigarra.up.pt/fcup/pt/noticias_geral.ver_noticia?p_nr=...
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 ...
- Result 2:**
 - Title:** "Curso Marketing Digital | Google"
 - Snippet:** "Notícia https://sigarra.up.pt/fcup/pt/noticias_geral.ver_noticia?p_nr=...
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 ...
- Result 3:**
 - Title:** "Google for Education"
 - Snippet:** "Notícia https://sigarra.up.pt/fpceup/pt/noticias_geral.ver_noticia?p_n...
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** ...
- Result 4:**
 - Title:** "Google for Education"
 - Snippet:** "Notícia https://sigarra.up.pt/fcup/pt/noticias_geral.ver_noticia?p_nr=...
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** ...

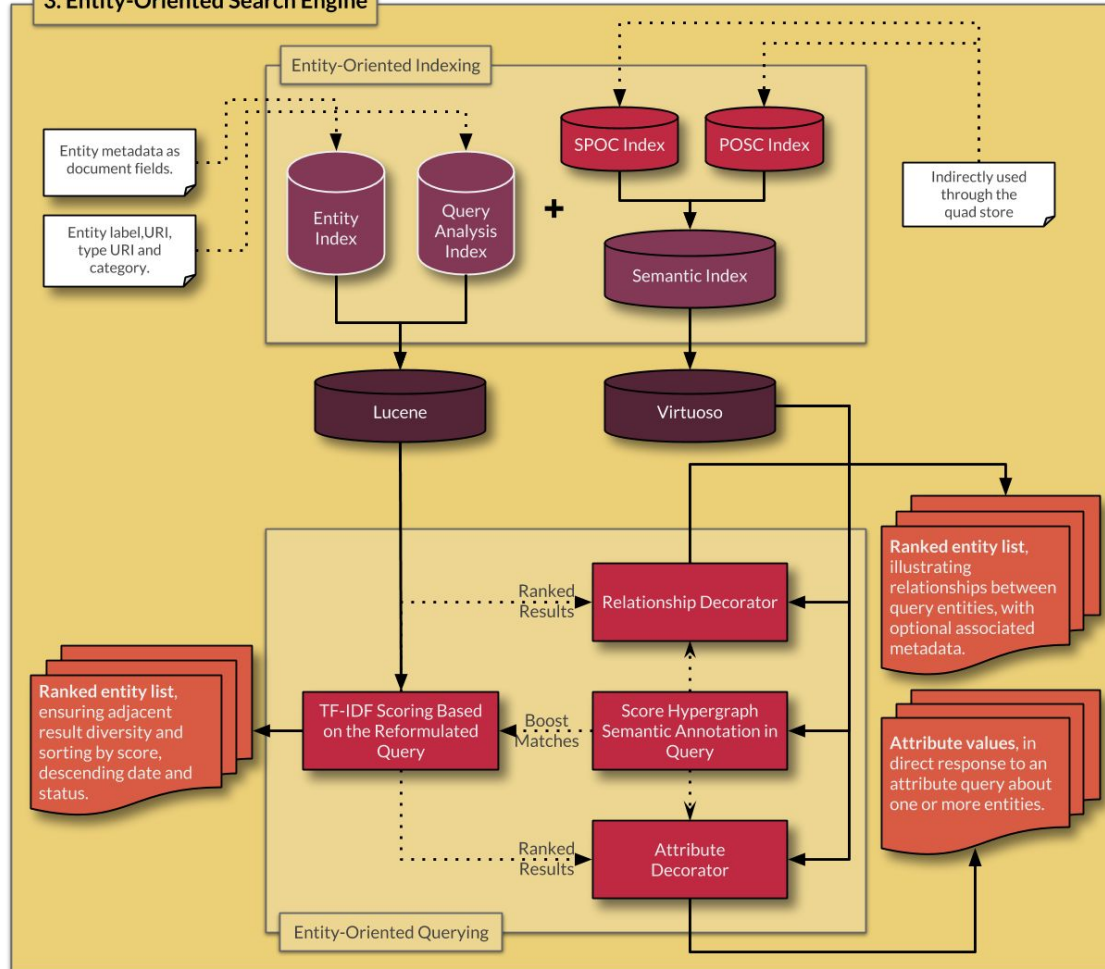
Keyword query. Anything that doesn't fit the previous four categories.

Search engine architecture

ANT components, from data collection to search.



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*

* Most relevant services.

Query understanding

What does it mean “to understand” and what is ANT’s approach?



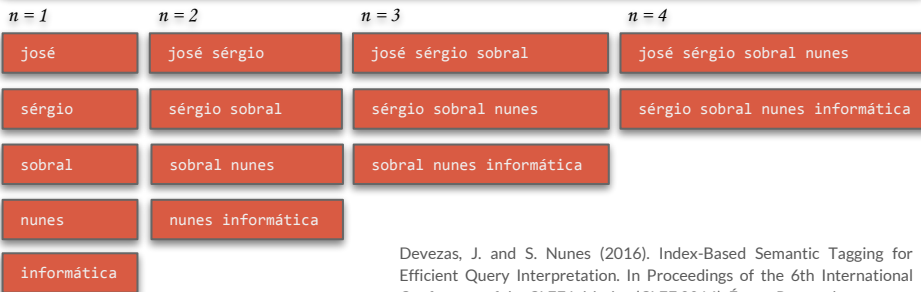
Query understanding

- Segment the query.
 - According to knowledge from the indexed corpus and from the knowledge base, which is the most likely division of the query?
- Annotate with semantic tags.
 - For each segment, assign the most likely semantic class (or none).
- Classify the query.
 - Based on the annotated segments, label the query with one of the five categories from Pound et al. (2010).
- In ANT, the query is seen as a sequence of keywords.
- But it could instead be considered as natural language:
 - Well-formed sentences;
 - Well-structured questions.
- A search engine's query processing methods should vary accordingly, with the style of query or, even, the user interface.

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



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.



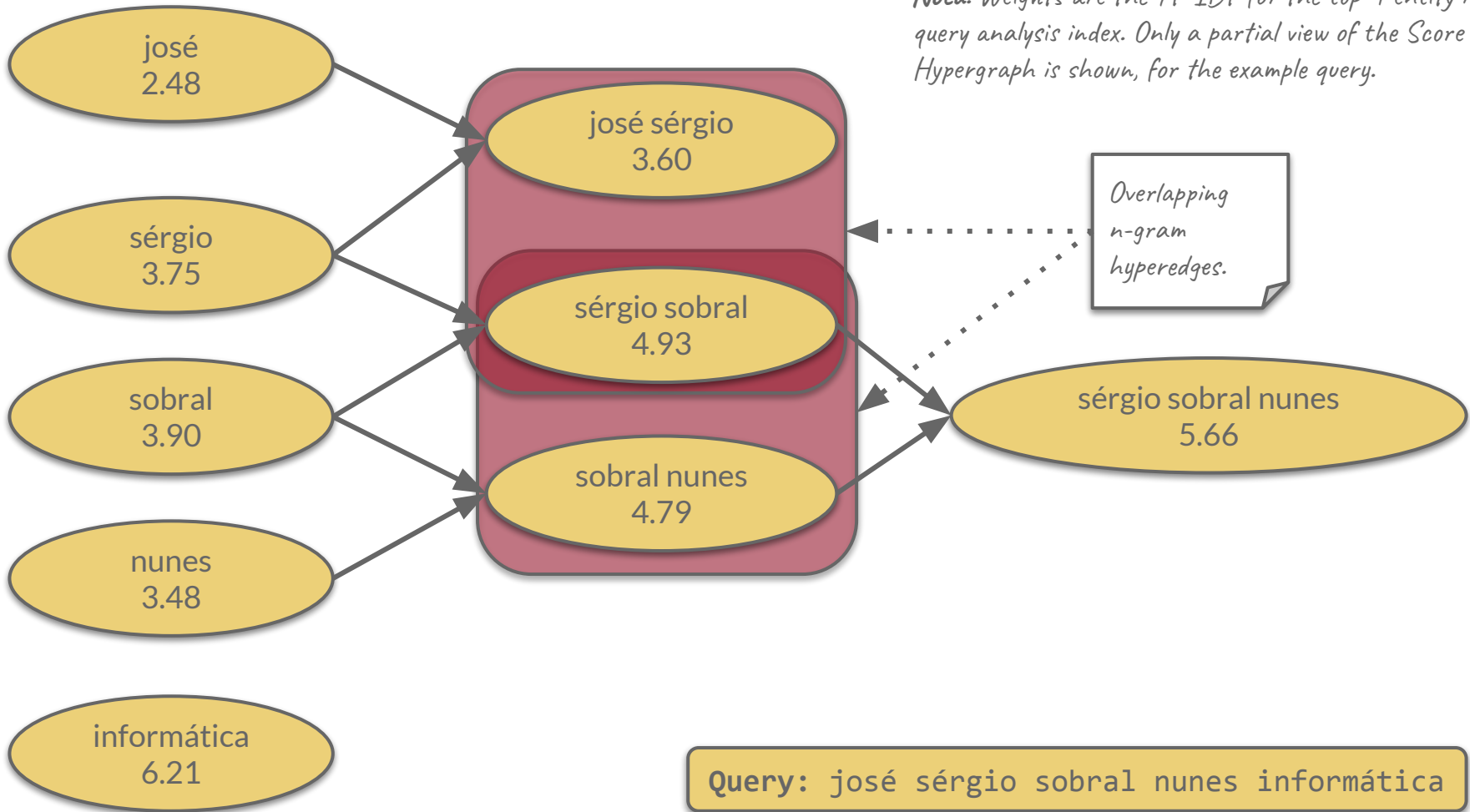
How does ANT understand queries?

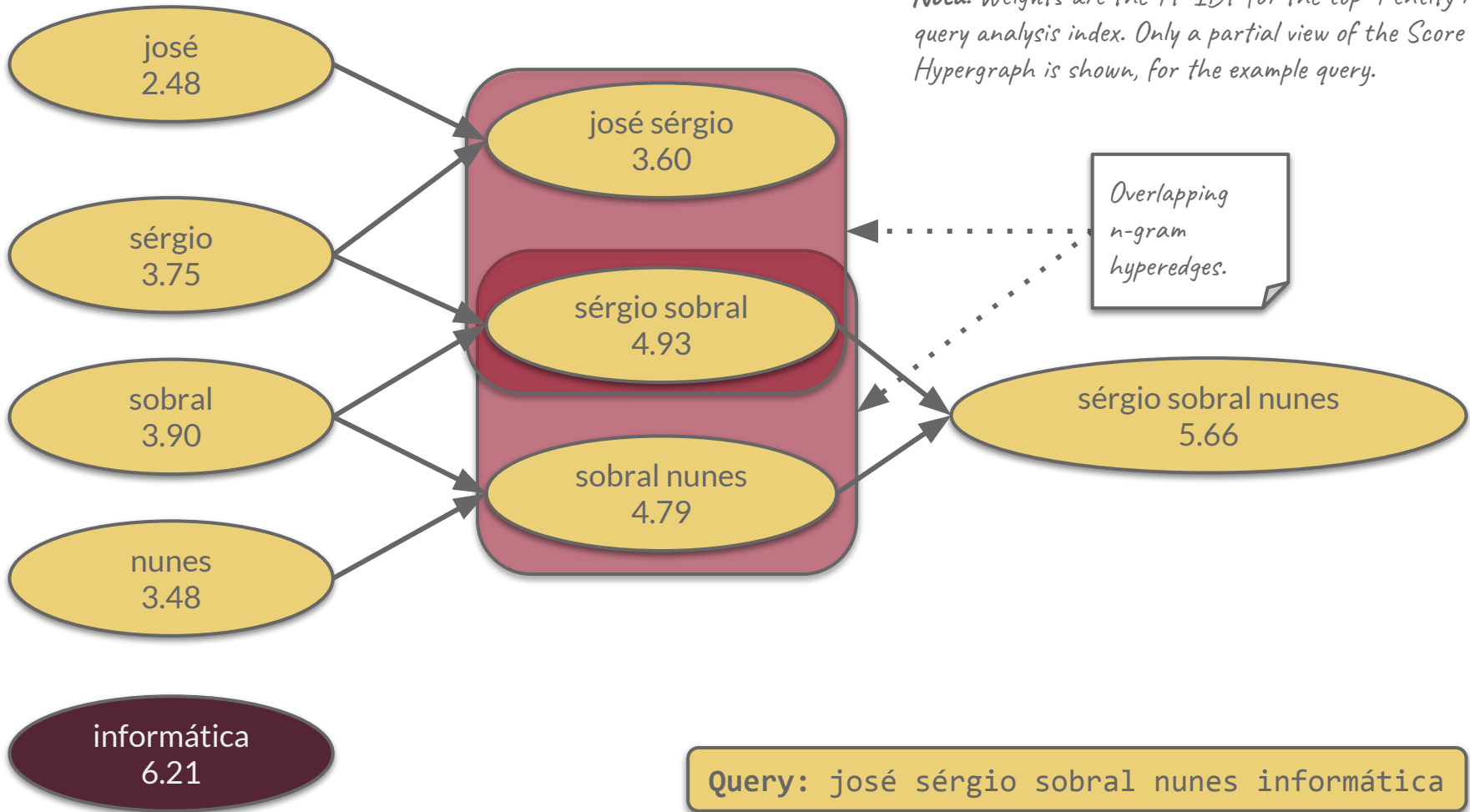
- In reality, we ended up using “**Score Hypergraph**” instead, which is a slight variation of the previous method, improving on performance and fixing some bugs.
- We used TF-IDF scores instead of probabilities.
- We created a dedicated index for query analysis, in order to search for entities matching n -grams, instead of querying the triple store.
- We used a hypergraph* of n -grams to solve query segment overlaps.

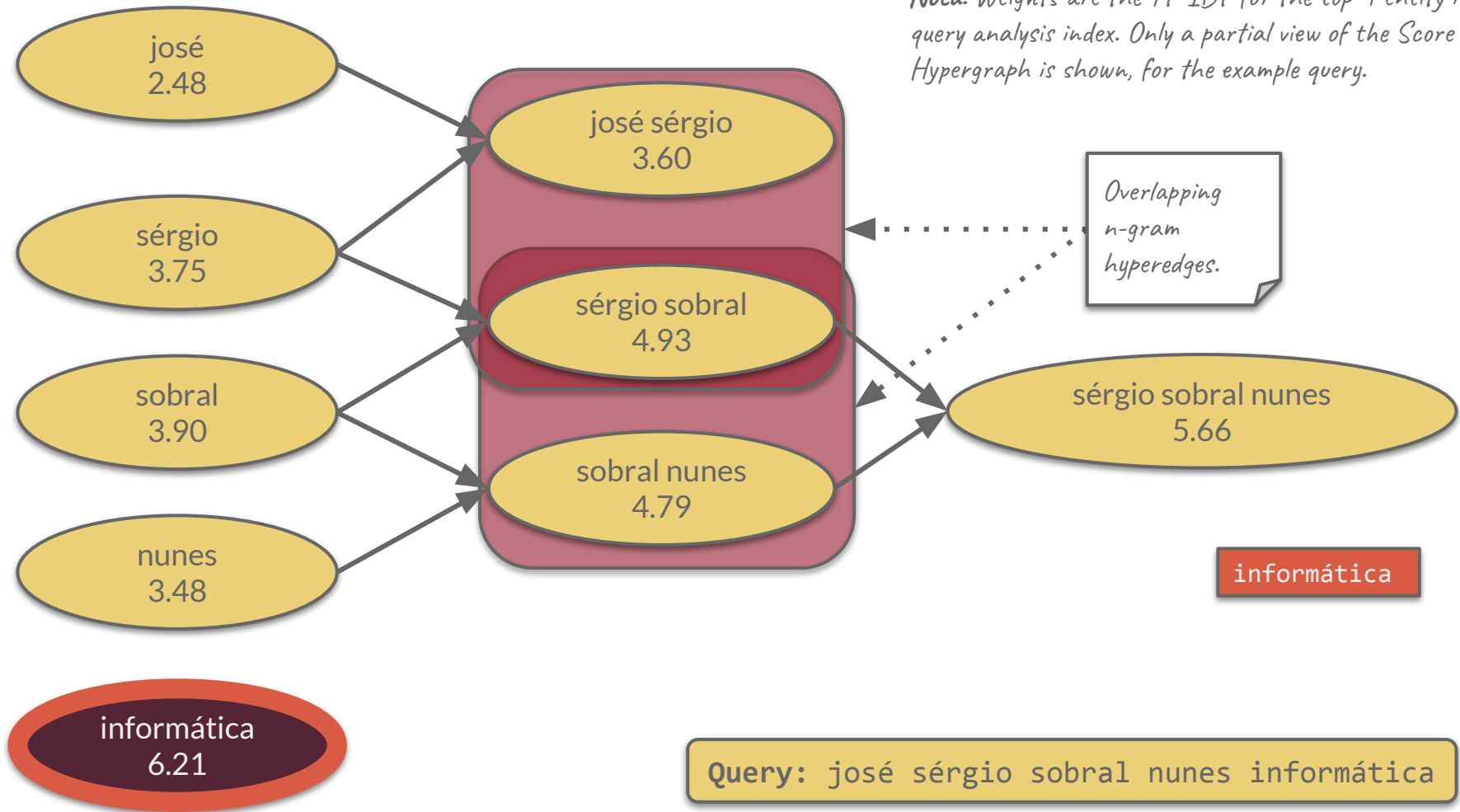
* A hypergraph is a generalization of a graph, where edges can have an arbitrary number of nodes.

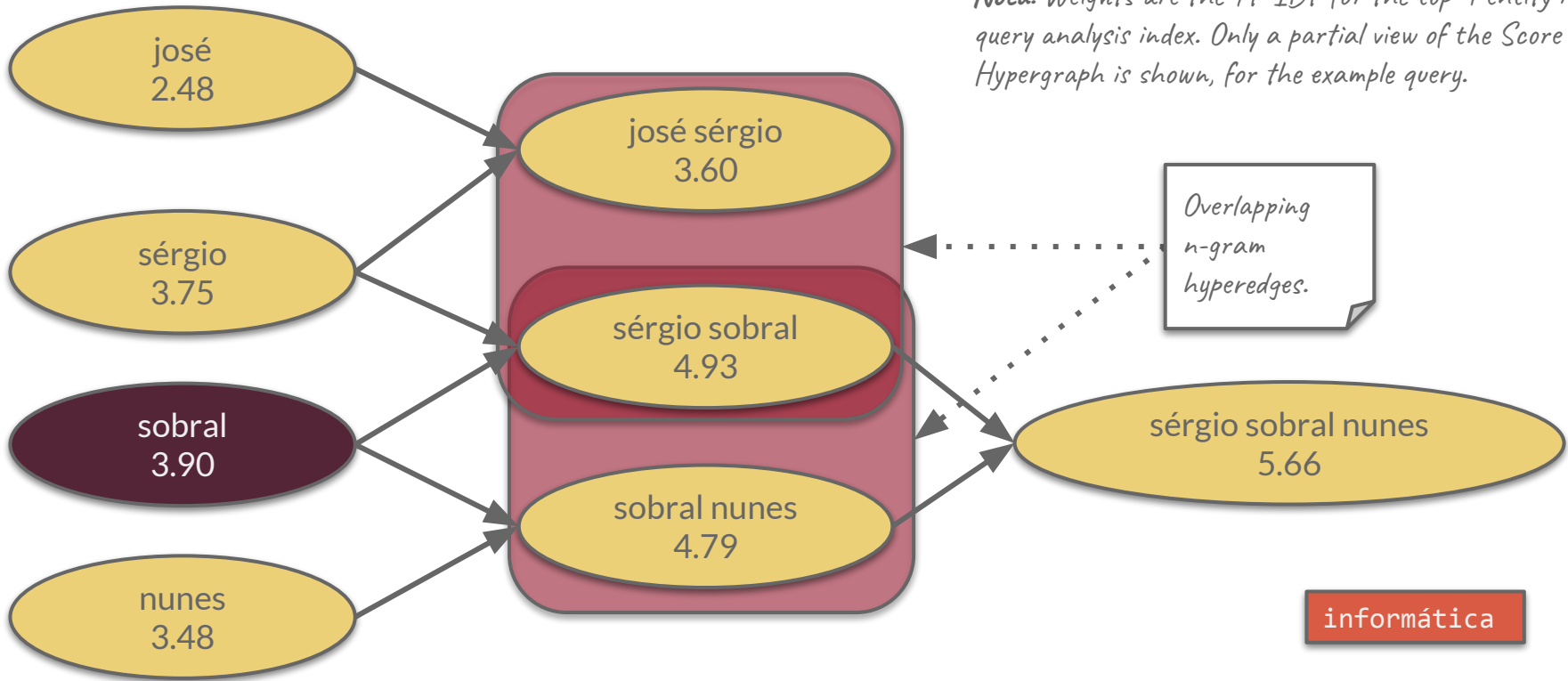
Score Hypergraph

Query segmentation and semantic tagging in ANT.



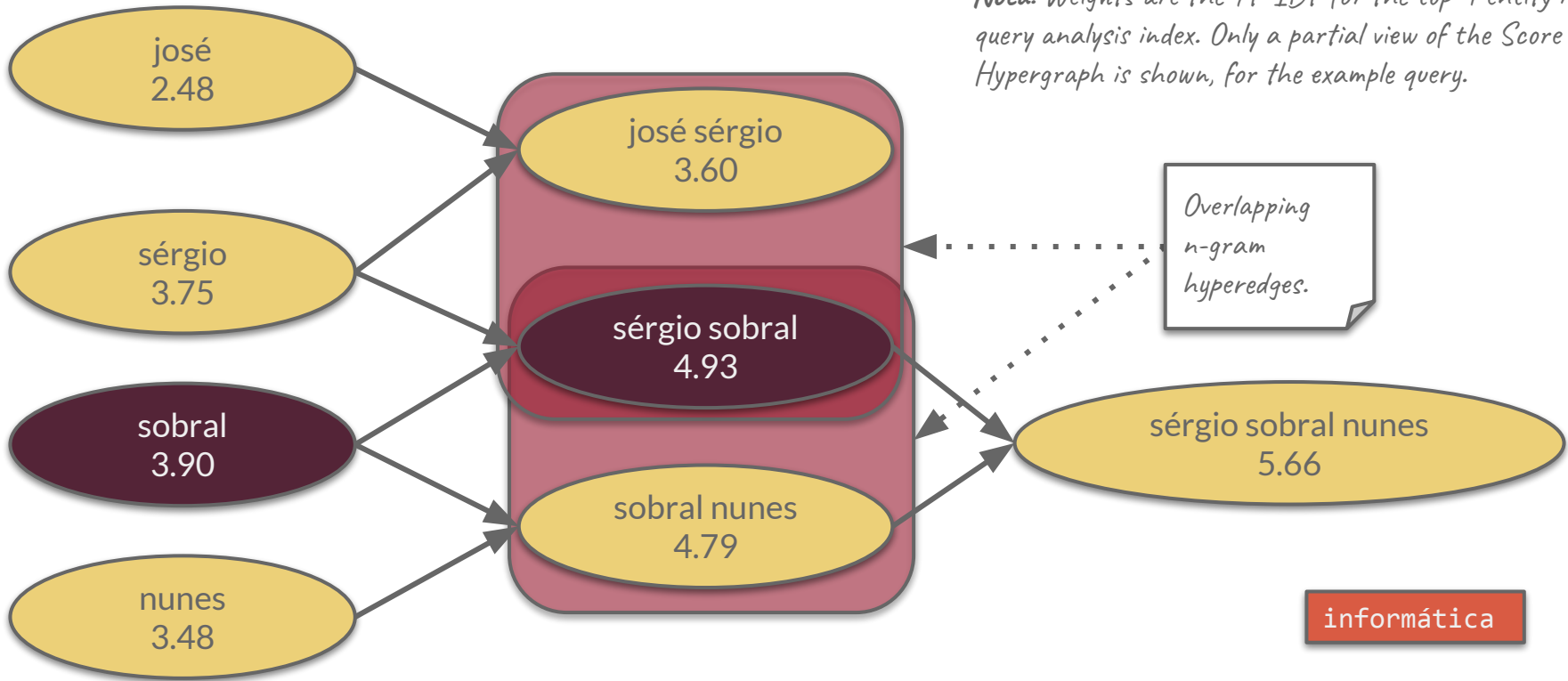






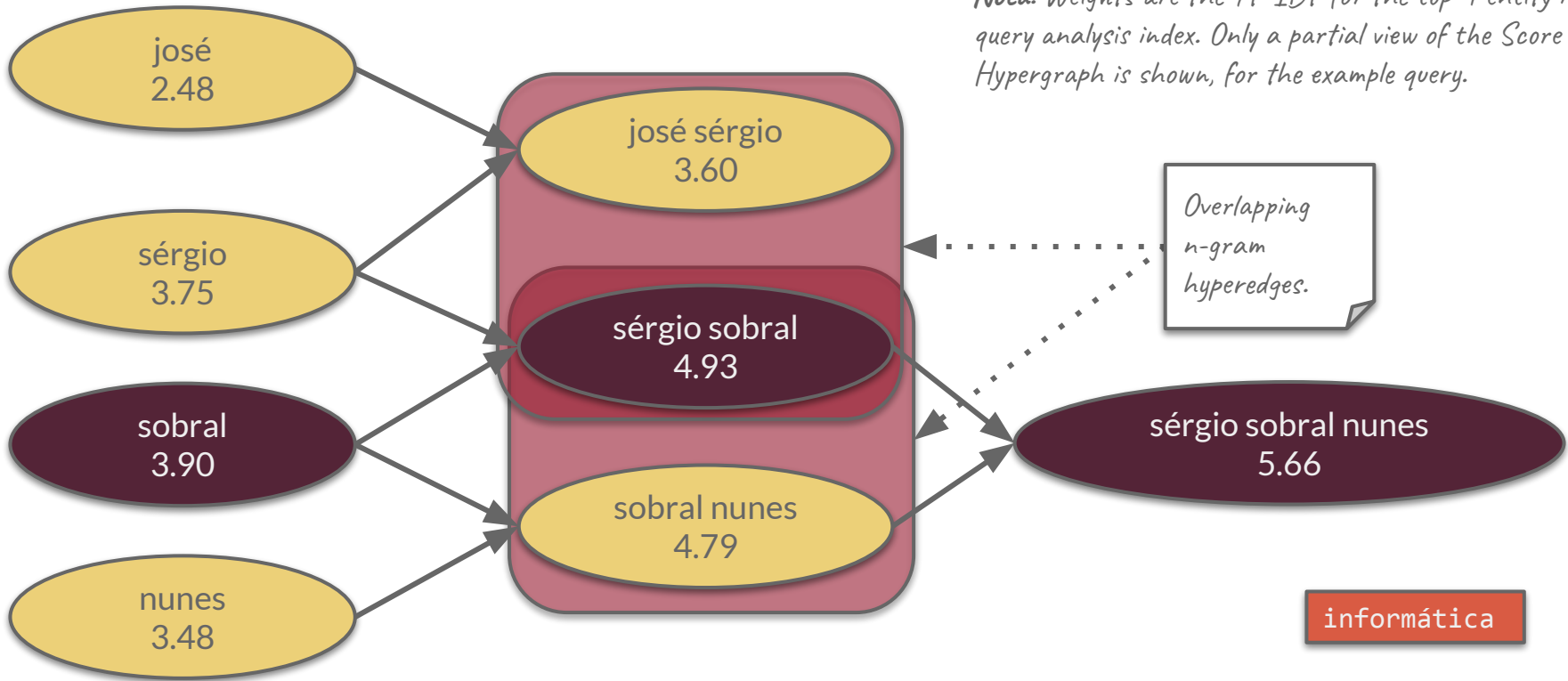
Nota: Weights are the TF-IDF for the top-1 entity in the query analysis index. Only a partial view of the Score Hypergraph is shown, for the example query.

Query: josé sérgio sobral Nunes informática



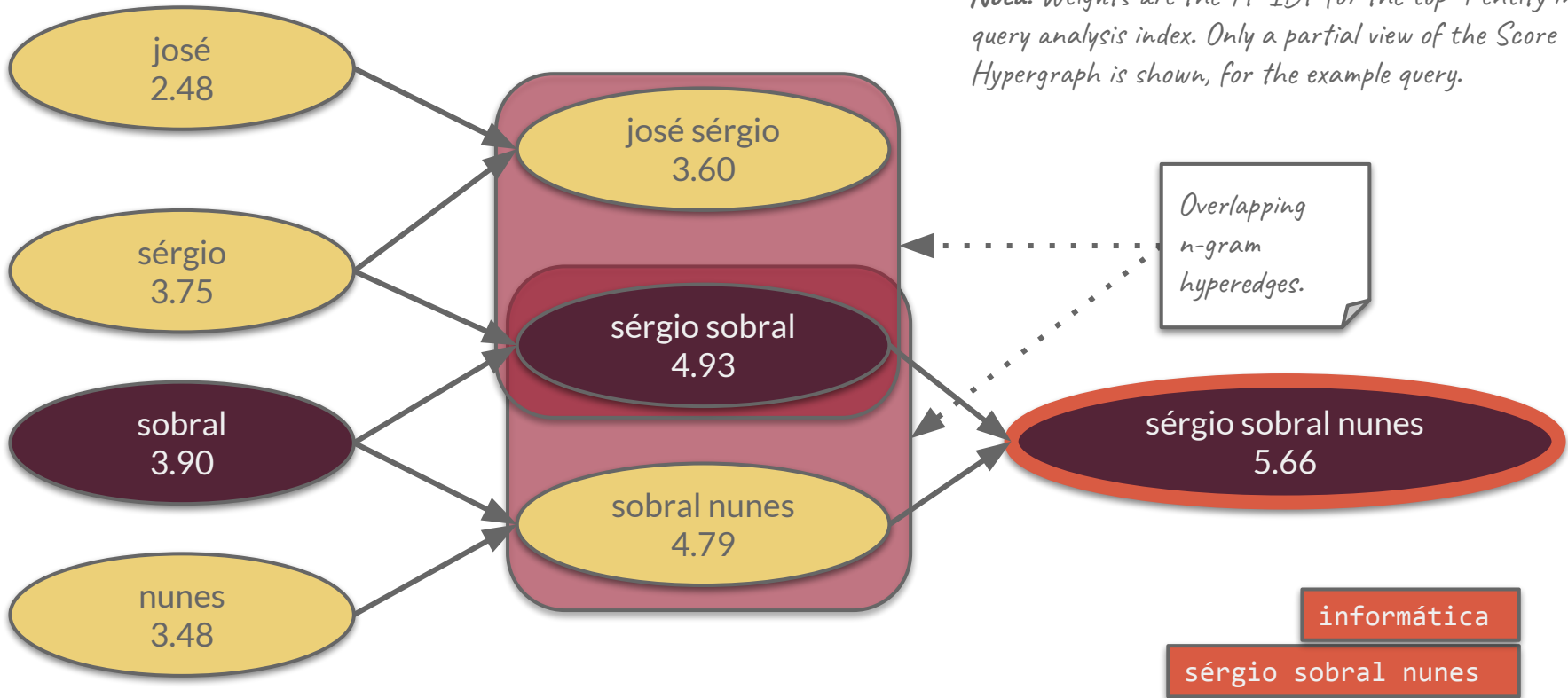
Nota: Weights are the TF-IDF for the top-1 entity in the query analysis index. Only a partial view of the Score Hypergraph is shown, for the example query.

Query: osé sérgio sobral nunes informática



Nota: Weights are the TF-IDF for the top-1 entity in the query analysis index. Only a partial view of the Score Hypergraph is shown, for the example query.

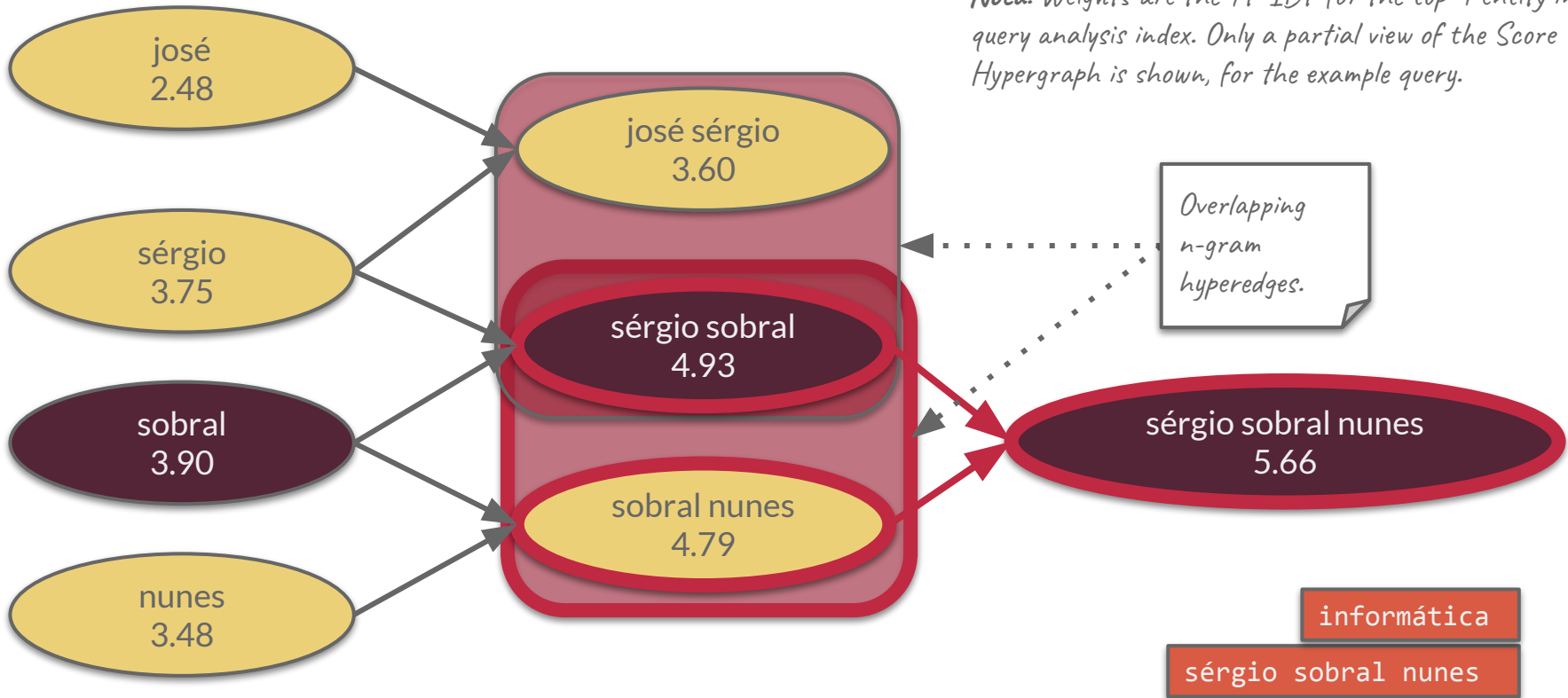
Query: josé sérgio sobral Nunes informática



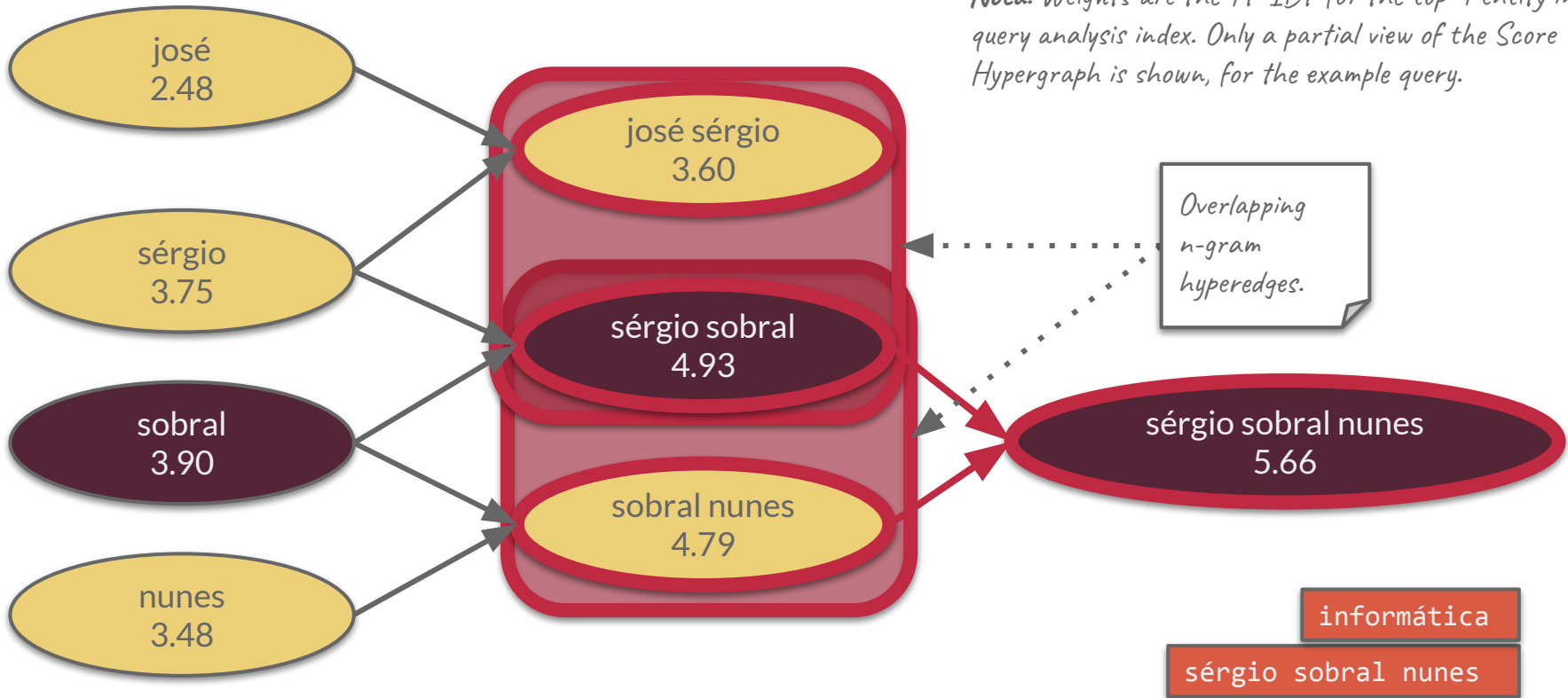
Nota: Weights are the TF-IDF for the top-1 entity in the query analysis index. Only a partial view of the Score Hypergraph is shown, for the example query.

*Overlapping
n-gram
hyperedges.*

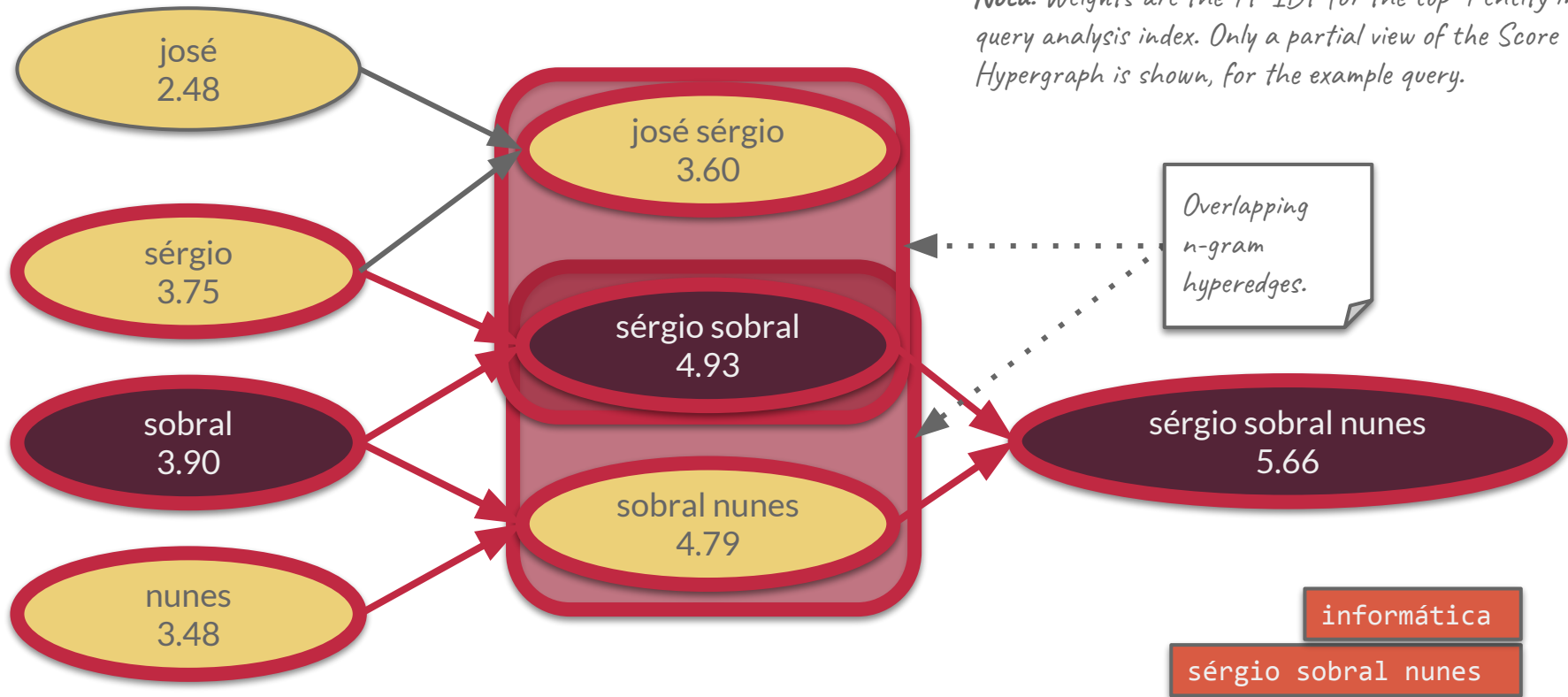
Query: josé sérgio sobral Nunes informática



Nota: Weights are the TF-IDF for the top-1 entity in the query analysis index. Only a partial view of the Score Hypergraph is shown, for the example query.



Nota: Weights are the TF-IDF for the top-1 entity in the query analysis index. Only a partial view of the Score Hypergraph is shown, for the example query.



Query: osé sérgio sobral nunes informática

josé
2.48

Nota: Weights are the TF-IDF for the top-1 entity in the query analysis index. Only a partial view of the Score Hypergraph is shown, for the example query.

informática

sérgio sobral nunes

Query: josé sérgio sobral nunes informática

josé
2.48

Nota: Weights are the TF-IDF for the top-1 entity in the query analysis index. Only a partial view of the Score Hypergraph is shown, for the example query.

informática

sérgio sobral nunes

Query: josé sérgio sobral nunes informática

josé
2.48

Nota: Weights are the TF-IDF for the top-1 entity in the query analysis index. Only a partial view of the Score Hypergraph is shown, for the example query.

informática

sérgio sobral nunes

josé

Query: josé sérgio sobral nunes informática

informática

sérgio sobral nunes

josé

Query: josé sérgio sobral nunes informática

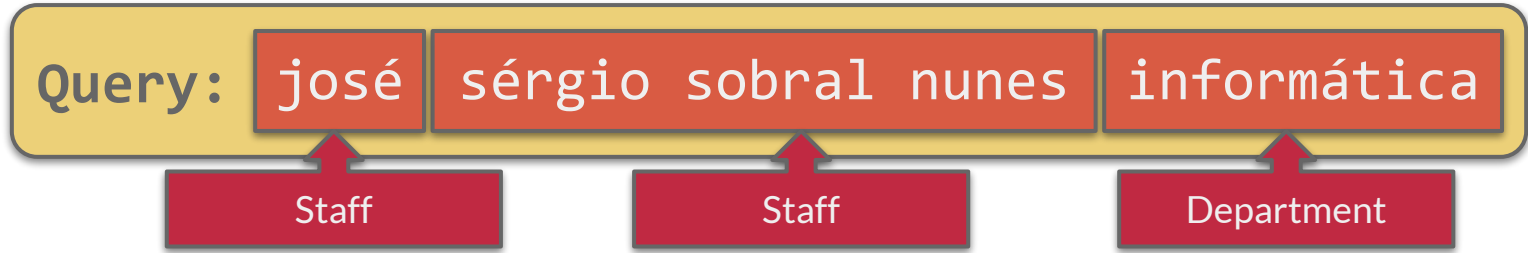
Query:

josé

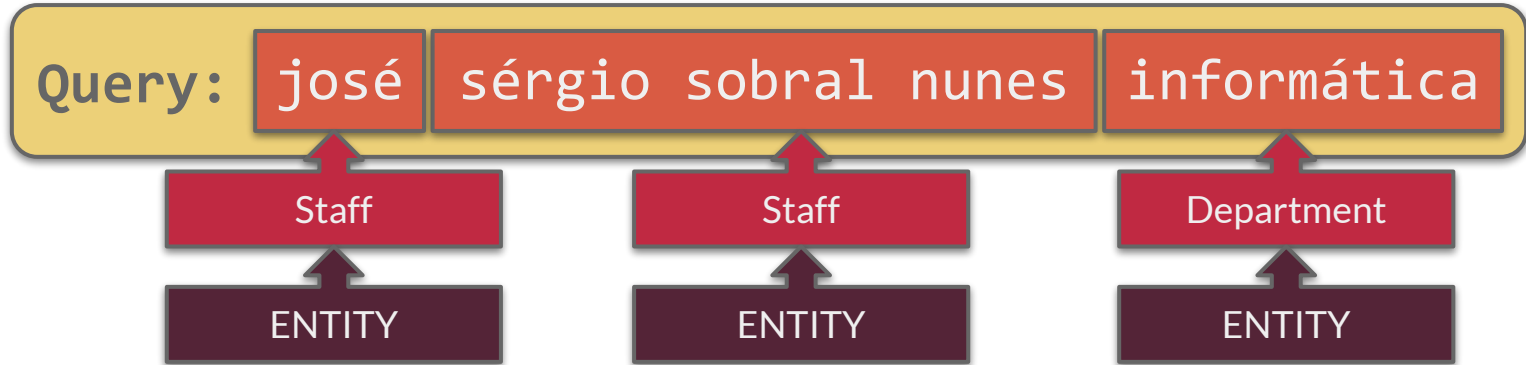
sérgio sobral nunes

informática

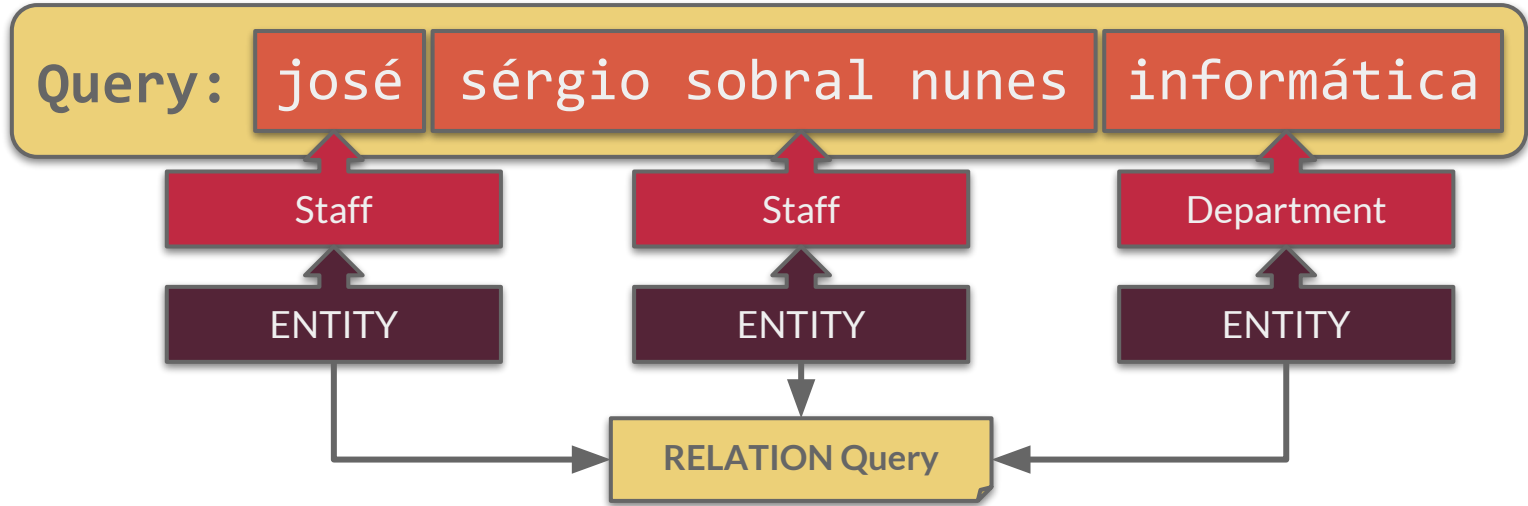
1. The query was segmented based on the n -grams with the highest-scoring entities.



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.

Conclusions

Final remarks, related projects, and a vision for the future.





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 has also served as 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.

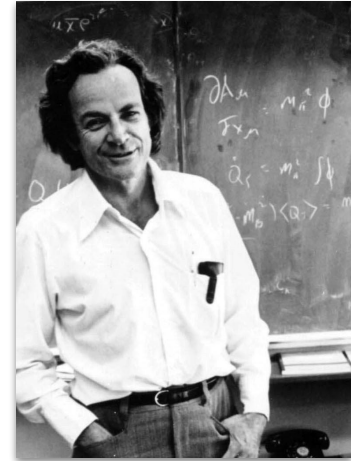


A vision for the future...

...can easily begin in a partially forgotten past.



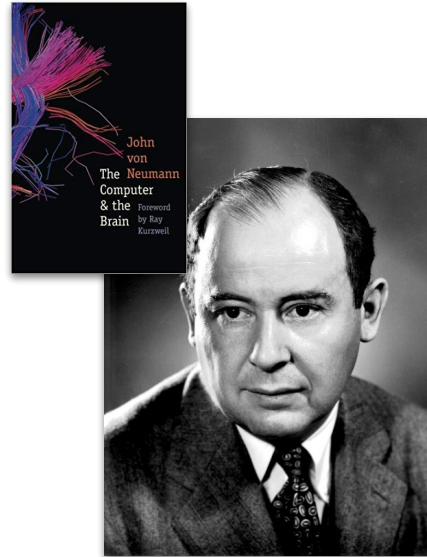
Von Neumann



Richard Feynman

A vision for the future...

...can easily begin in a partially forgotten past.



Von Neumann

Prior to his death, in 1957, John Von Neumann was preparing his classes for the Yale Silliman lectures. The idea was to compare the computer and the brain, studying both of their organs (neuron/transistor).

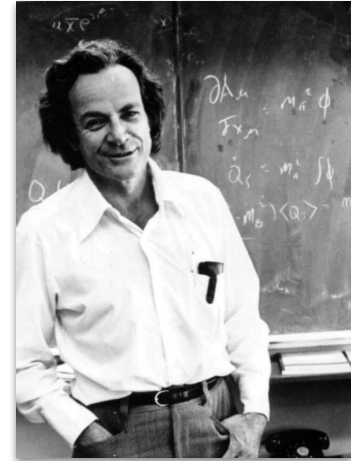
In the present, Cognitive Science already fits the role of tying together Neuroscience and Artificial Intelligence. Nevertheless, there are still unexplored paths at the intersections. There are answers and ideas that have been lost in the past!

A vision for the future...

...can easily begin in a partially forgotten past.

Theoretical physicist, with relevant contributions in Quantum Mechanics.

During World War II, he administered the group of human computers in the theoretical division of the Los Alamos Laboratory. Human computing is nothing else than manually solving a set of long and complex calculations, in group and strictly following pre-established rules.



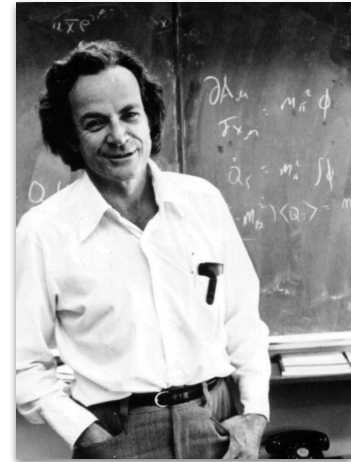
Richard Feynman

A vision for the future...

...can easily begin in a partially forgotten past.

Maybe that's why he collaborated with Stanley Frankel and Nicholas Metropolis in developing a system for using IBM punch cards for computation.

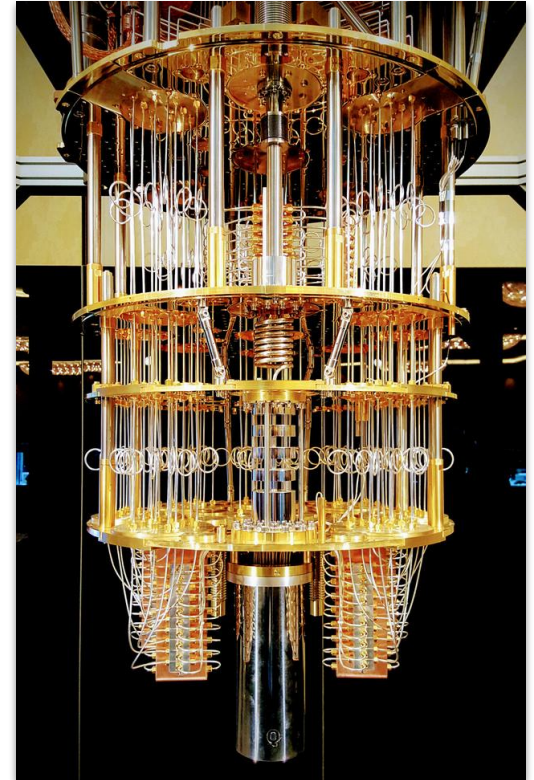
He also invented a new method for computing logarithms that he later used in the Connection Machine, developed by Daniel Hillis as an alternative to the Von Neumann architecture.



Richard Feynman

A vision for the future...

...towards quantum search engines?

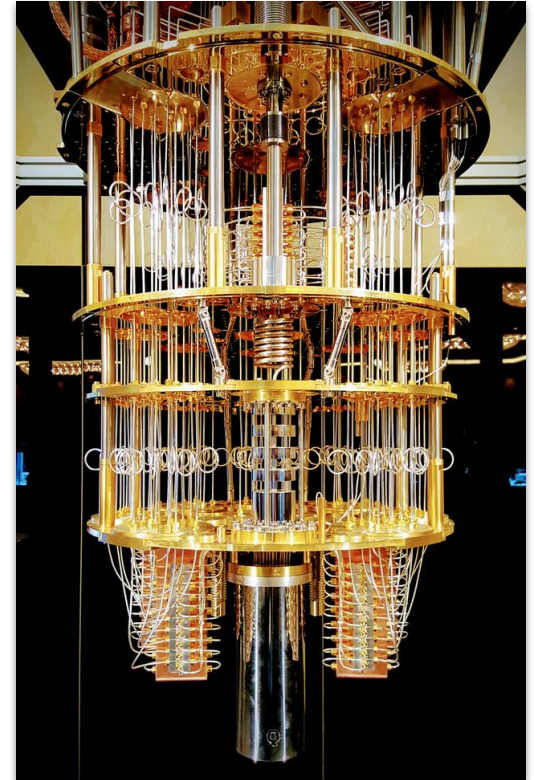


A vision for the future...

...with many unanswered questions.

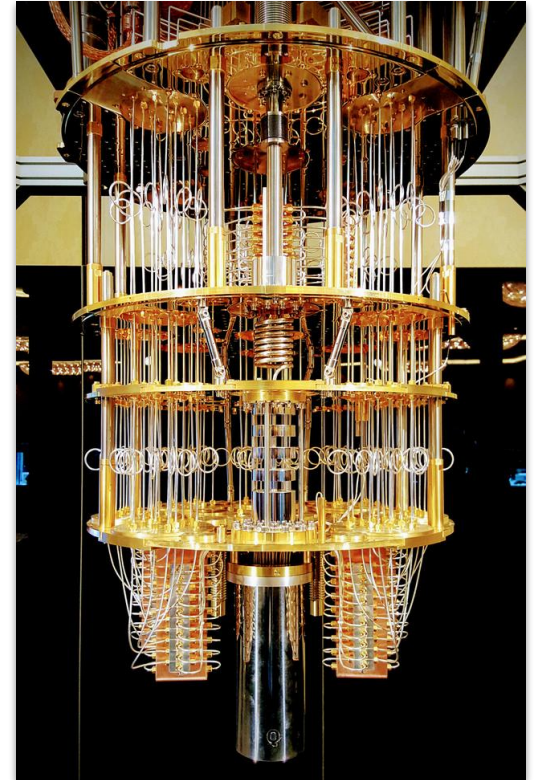
“Major Quantum Computing Advance Made Obsolete by Teenager”

<https://www.quantamagazine.org/teenager-finds-classical-alternative-to-quantum-recommendation-algorithm-20180731/>



A vision for the future...

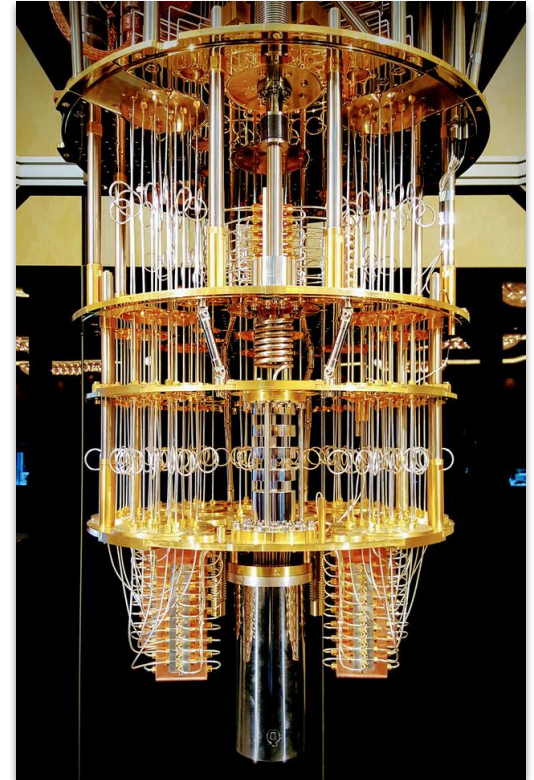
...where uncertainty and probabilities trades places with binary.



A vision for the future...

...where uncertainty-based search reinforces its relevance.

At FEUP InfoLab, we're working with classical computers, while exploring ideas that might have quantum applications.



A vision for the future...

...that is happier, more creative and increasingly exploratory.

The time is now, to combine the tools that science has already provided. The future is hybrid!



You mean... Electric?

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.

