

Graph-Based Entity-Oriented Search

José Devezas <joseluisdevezas@gmail.com>

Thesis supervisor: Sérgio Nunes

INESC TEC and Faculty of Engineering, University of Porto

Thesis submitted to Faculty of Engineering of the University of Porto for the Doctor Degree in Computer Science within the Joint Doctoral Program in Computer Science of the Universities of Minho, Aveiro and Porto.

Porto, Portugal – January 26, 2021

What to watch

Services

Movies

Sci-fi

1980s

Adventure

Horror

Comedy

YouTube

Drama

Google Play

Superhero

Love

Thriller

Disney

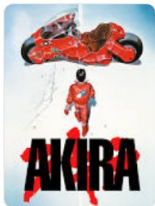
v



Blade Runner



1984



Akira



E.T. the Extra-Terrestrial



Back to the Future



The Terminator



Predator



Aliens



More movies

Feedback

About 20,100,000 results (1.16 seconds)

collider.com › Movie › The Best 80s Sci-Fi Movies

The Best 80s Sci-Fi Movies - Collider

Mar 8, 2019 — **The Empire Strikes Back** (1980) No Star Wars movie since has come close to delivering such a complete and well-told story as Empire Strikes Back. **Flash Gordon** (1980) **Scanners** (1981) **The Road Warrior** (1981) **Time Bandits** (1981) **Escape From New York** (1981) **E.T.** (1982) **Star Trek II: The Wrath of Khan** (1982)

en.wikipedia.org › wiki › List_of_science_fiction_films...

List of science fiction films of the 1980s - Wikipedia

A list of **science fiction** films released in the **1980s**. These films include core elements of **science fiction**, but can cross into other genres. They have been ...

What to watch

Services

Movies

Sci-fi

1980s

Adventure

Horror

Comedy

YouTube

Drama

Google Play

Superhero

Love

Thriller

Disney

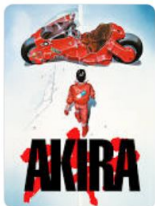
v



Blade Runner



1984



Akira



E.T. the Extra-Terrestrial



Back to the Future



The Terminator



Predator



Aliens

Entities



More movies

Feedback

About 20,100,000 results (1.16 seconds)

collider.com › Movie › The Best 80s Sci-Fi Movies

The Best 80s Sci-Fi Movies - Collider

Mar 8, 2019 — **The Empire Strikes Back** (1980) No Star Wars movie since has come close to delivering such a complete and well-told story as Empire Strikes Back. **Flash Gordon** (1980) **Scanners** (1981) **The Road Warrior** (1981) **Time Bandits** (1981) **Escape From New York** (1981) **E.T.** (1982) **Star Trek II: The Wrath of Khan** (1982)

en.wikipedia.org › wiki › List_of_science_fiction_films...

List of science fiction films of the 1980s - Wikipedia

A list of **science fiction** films released in the **1980s**. These films include core elements of **science fiction**, but can cross into other genres. They have been ...

Documents

Type

Attributes

Services

Movies

Sci-fi

1980s

Adventure

Horror

Comedy

YouTube

Drama

Google Play

Superhero

Love

Thriller

Disney

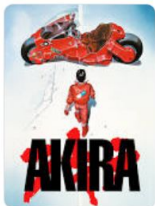
v



Blade Runner



1984



Akira

E.T. the Extra-
Terrestrial

Back to the Future



The Terminator



Predator



Aliens

Entities



More movies

Feedback

About 20,100,000 results (1.16 seconds)

collider.com › Movie › The Best 80s Sci-Fi Movies

The Best 80s Sci-Fi Movies - Collider

Mar 8, 2019 — **The Empire Strikes Back** (1980) No Star Wars movie since has come close to delivering such **Entity** story as Empire Strikes Back. Flash Gordon (1980) Scanners (1981) The Road Warrior (1981) Time Bandits (1981) Escape From New York (1981) E.T. (1982) Star Trek II: The Wrath of Khan (1982)

en.wikipedia.org › wiki › List_of_science_fiction_films...

List of science fiction films of the 1980s - Wikipedia

A list of **science fiction** films released in the **1980s**. These **films** include core elements of **science fiction**, but can cross into other genres. They have been ...

Documents

Entity-oriented search is the search paradigm of organizing and accessing information centered around entities, and their attributes and relationships.

– Krisztian Balog, 2018

- Entities and their relations
- Documents mentioning entities

- Knowledge bases

- Corpora

- Triplestores
- Inverted indexes

- Triplestores



Opportunity for a
unified framework

- Inverted indexes

- Structured data and queries



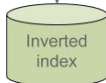
Opportunity for a
unified framework

- Unstructured data and queries

Query: entertainers that are friends with astronauts who walked on the moon

"After the act, Kevin Foster went down to the audience, to hug his friend, Neil Armstrong, who had been sitting in the crowd since the beginning of the show."

Neil Armstrong	:isA	Astronaut
Neil Armstrong	:walkedOn	Moon
Buzz Aldrin	:isA	Astronaut
Buzz Aldrin	:walkedOn	Moon
Kevin Foster	:isA	Entertainer



Unstructured query:
entertainer friend
astronaut walk moon

Structured query:
SELECT ?e
WHERE {
 ?e a :Entertainer .
 ?e :friend ?a .
 ?a a :Astronaut .
 ?a :walkedOn :Moon .
}

Ranked documents

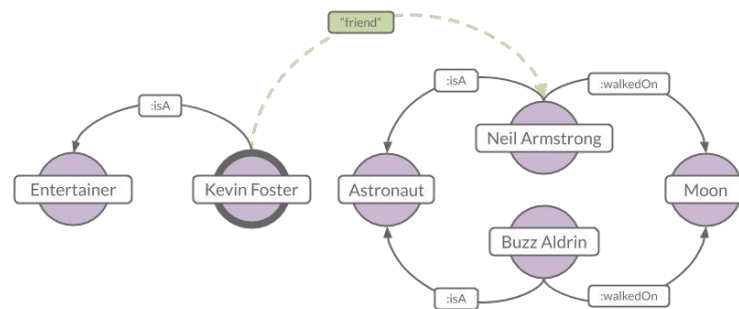
Documents containing the terms 'astronaut', 'walk' and 'moon' should rank better than our "niche" document of interest matching only 'friend'.

Unranked entities
(boolean model)



NER

Ranked entities
(rank fusion / score weight
combination)



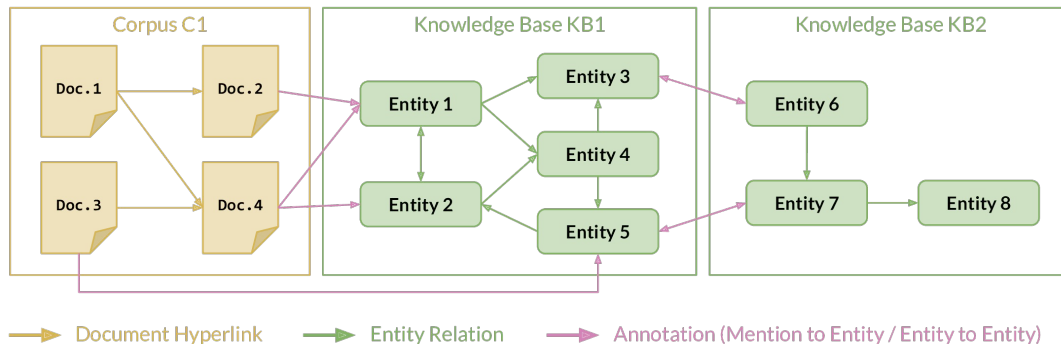
Motivation

- Separate representation models requires integration
- Opportunity for a joint model
 - Reach a wider range of answers
 - Generalize retrieval tasks

A unified model for entity-oriented search

Combined data

- Text
- Entities
- Relations



A unified model for entity-oriented search

Retrieval tasks

- Ad hoc document retrieval
- Ad hoc entity retrieval
- Related entity finding
- Entity list completion

A unified model for entity-oriented search

Retrieval tasks

- Ad hoc document retrieval
- Ad hoc entity retrieval
- Related entity finding
- Entity list completion

Input

Keyword Query: croft bendersky

Document mentioning the College of Information and Computer Sciences and Hypergraph entities, since W. Bruce Croft is dean of the College of Information and Computer Sciences and Hypergraph is one of the topics covered by Michael Bendersky in his thesis.

Output

Doc:
342

Doc:
13

Doc:
671

Faculty member and former Dean in the College of Information and Computer Sciences.

Recent Ph.D. Graduates:

Michael Bendersky
Van Dang

A unified model for entity-oriented search

Retrieval tasks

- Ad hoc document retrieval
- Ad hoc entity retrieval
- Related entity finding
- Entity list completion

Input

Keyword Query: croft bendersky

Output

Entity: [Person] W. Bruce Croft

Entity: [Person] Michael Bendersky

A unified model for entity-oriented search

Retrieval tasks

- Ad hoc document retrieval
- Ad hoc entity retrieval
- Related entity finding
- Entity list completion

Input

Entity: [Person] Michael Bendersky

Type: [ScholarlyArticle]

Relation: [creator]

Output

Entity: [ScholarlyArticle] Discovering key concepts in verbose queries

Entity: [ScholarlyArticle] Modeling higher-order term dependencies in information retrieval using query hypergraphs

A unified model for entity-oriented search

Retrieval tasks

- Ad hoc document retrieval
- Ad hoc entity retrieval
- Related entity finding
- Entity list completion

Input

Entity: [Person] Michael Bendersky

Type: [ScholarlyArticle]

Relation: [creator]

Example 1: [ScholarlyArticle] Information retrieval with query hypergraphs

Output

Entity: [ScholarlyArticle] Modeling higher-order term dependencies in information retrieval using query hypergraphs

Entity: [ScholarlyArticle] Discovering key concepts in verbose queries

This is more similar to the example, so we moved it up.

THESIS STATEMENT

Graphs can be used to jointly index corpora and knowledge bases, supporting retrieval for multiple entity-oriented search tasks.

Main objectives

- Joint representation of terms, entities, and their relations
- Universal ranking function for multiple entity-oriented search tasks
- Improved retrieval effectiveness through the unification of information sources



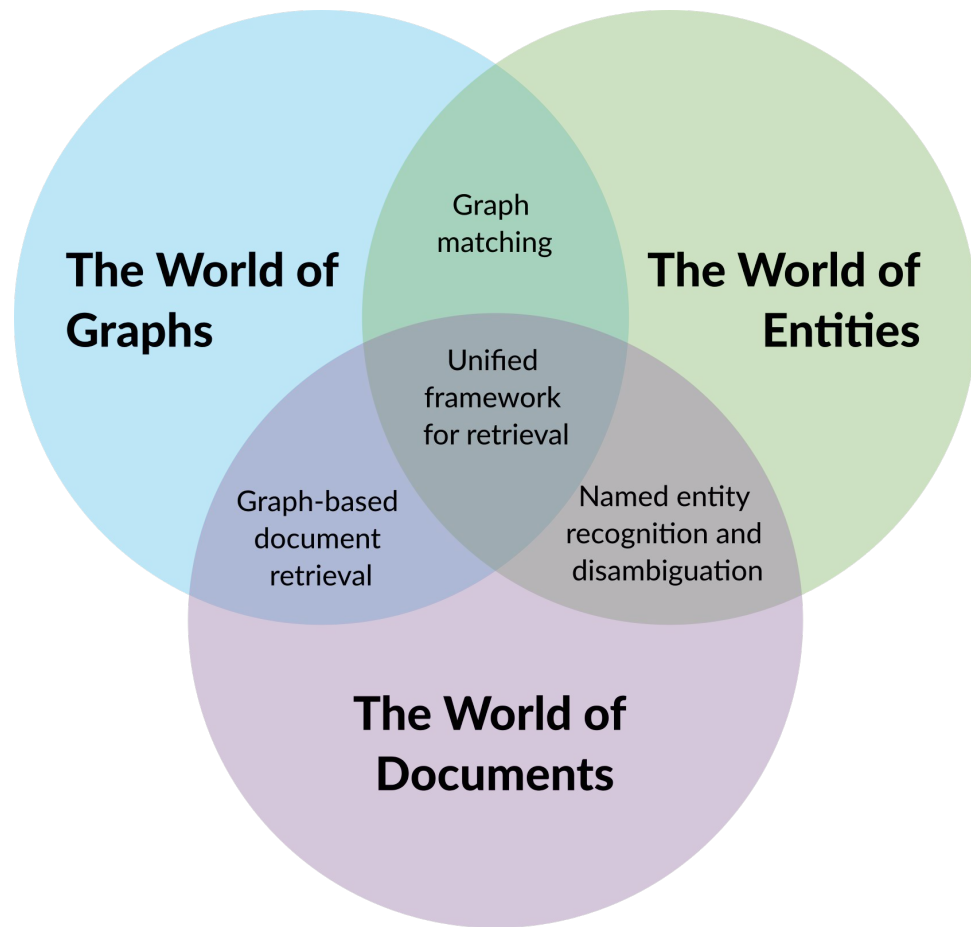
State of the art

A breadth-first search for intersecting concepts in the worlds of documents, entities, and graphs.

Three axes covered

Research around a “common denominator”:

- Documents as graphs of words
- Entity relations as knowledge graphs
- Generic graph-based models and applications



Key references

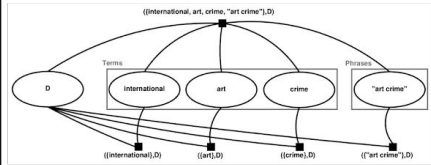
2007

"73-87% of all queries contain entities"

Concordance-Based
Entity-Oriented Search
Bautin, M. and Sklena, S.

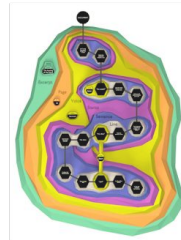
2010

2012



Modeling Higher-order Term
Dependencies in Information
Retrieval Using Query Hypergraphs
Bendersky, M. and Croft, W. B.

2013

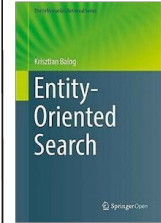


Index for Efficient Semantic
Full-text Search
Bast, H. and Buchhold, B.

2017

It's more than just
overlap: Text As Graph
Dekker, R. H. and Birnbaum, D.

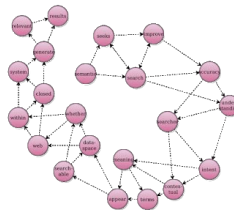
2018



Entity-Oriented Search
Balog, K.



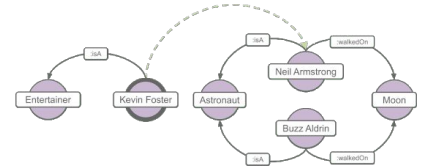
Ad-hoc object retrieval
in the web of data
Pound, J., Mika, P. and Zaragoza, H.



Query: entertainers that are friends with astronauts who walked on the moon

"After the act, Kevin Foster went down to the audience, to hug his friend, Neil Armstrong, who had been sitting in the audience since the beginning of the show."	Neil Armstrong :isA Astronaut	Astronaut
	Neil Armstrong :walkedOn Moon	Moon
	Buzz Aldrin :isA Astronaut	Astronaut
	Buzz Aldrin :walkedOn Moon	Moon
	Kevin Foster :isA Entertainer	Entertainer

Graph-of-word and TW-IDF: new
approach to ad hoc IR
Rousseau, F. and Vazirgiannis, M.



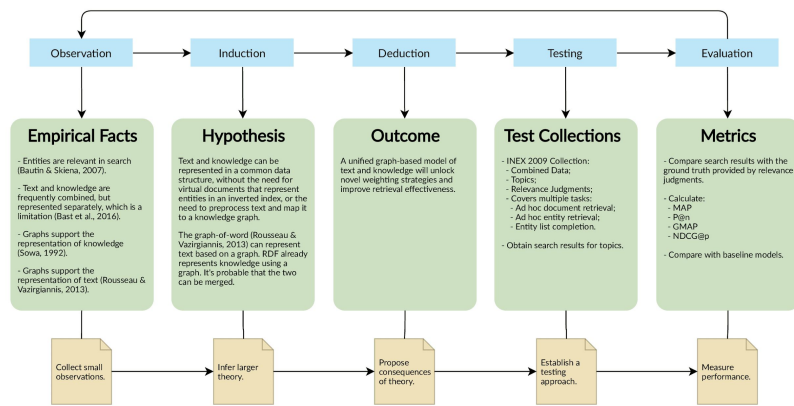
II

Materials and methods

Empirical research supported on test
collections and software

Research methodology

Empirical research



Systematic documentation

INEX 2009 Wikipedia Collection

Global	
Source	http://www.mim.uni-saarland.de/data/documentcollection.html#wikipedia
Paper	http://dx.doi.org/10.1007/978-3-642-01031-7_103
Date	October 8, 2008
Size	5.5 GB compressed; 50.7 GB uncompressed

Statistics	
Documents	2,666,190
Entities	101,917,424 XML elements per document
Topics	115 for 2009; 107 for 2010
Assessments	50,725 for 2009; 36,031 for 2010

Description

Starting in 2009, INEX uses a new document set which has been converted into XML, using both general tags (for text and term), typographical tags (for bold, italic, etc.) and semantic annotations from the 2005 w4C2 version. The collection was created from the October 8, 2008 snapshot of the Wikipedia database. It contains a total of 5.5 GB of text (50.7 GB of XML). There are (excluding white-space).

Conference

INEX 2009: The 10th International Conference on Web Intelligence

Year

2007

To Do

Is it possible to justify the relevance of entities in queries (i.e., why is entity-oriented search relevant)?

Review extended version

Review extended version.

Hypergraph-of-Entity

ID	Experiment 2
Start Date	2017-10-24 16:38
End Date	Ongoing
Why do it?	The graph-of-entity exploded in number of edges. Using hyperedges might enable the aggregation of multiple edges in a single edge, reducing dimension and allowing for the planned experiments.
Main strengths	Indexing can be done in about 3 minutes for an in-memory version of the graph.

Table of Contents

- Concordance-Based Entity-Oriented Search
- Introduction
- Related Work
- Entities in Web Queries
- Concordance-Based Entity Search
- Evaluation
- Conclusions and Future Work

Table of Contents

- Hypergraph-of-Entity
- Challenges
- Dependencies
- Visuals
- Traces
- Evaluation
- Test Runs
- Start 2018
- End 2019 (see OCS 2019 for corrections)
- OCS 2019
- TBD
- Comparing with Graph-of-Entity
- Representation Ranker
- ECR 2020: Completely Indexing INEX 2009
- Research Log

Collections

Bautin and Skiena present what they consider to be the "first-in-literature" implementation of an entity search engine. Their contribution is quite rich, in the sense that they cover multiple facets of entity-oriented search. In particular, they found that nearly 87% queries contain entities, by analyzing the AOL dataset with 30 million web search queries. They also proposed a concordance-based model for entity representation, along with an adaptation of Lucene's TF-IDF scoring scheme, where each document (a concordance) is a concatenation of all sentences containing a given entity, optionally for a given period of time (e.g., month). Thus, they also propose a time-dependent scoring function, modeling user interest in an entity as a function of time, optimizing parameters based on the frequency of entities in the AOL query log. Finally, they propose a method for evaluating an entity search engine, by comparing the results list with the equivalent list obtained through a juxtaposition score. The juxtaposition score measures the upper bound of the probability of two entities occurring in the same sentence under the assumption of independence. By obtaining the results list from Lucene and the results list based on top related entities according to juxtaposition, the lists are then compared using the $A_{w,d}$ distance from Fagin et al., showing the best results for phrase queries with the stop parameter (i.e., word-based edit distance) equal to the number of query terms.

Experiments

is, this is still too inefficient.

on-WALK_LENGTH and WALK_REPEATS (compare ranking lists)

quantification of concordance.

is integration into the models.

test several functions until we find a basic one with a good

define thresholds to prune (e.g., bottom 10%).

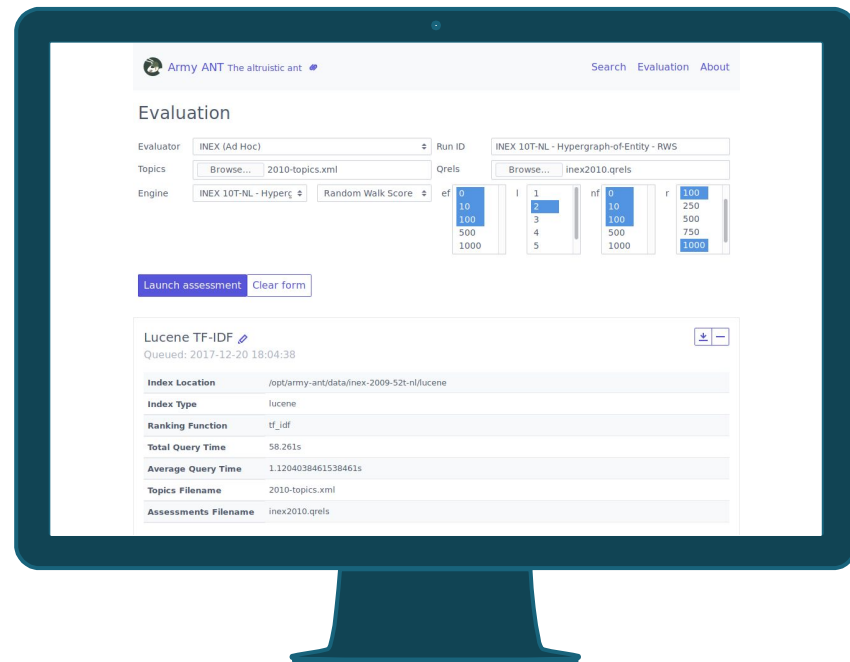
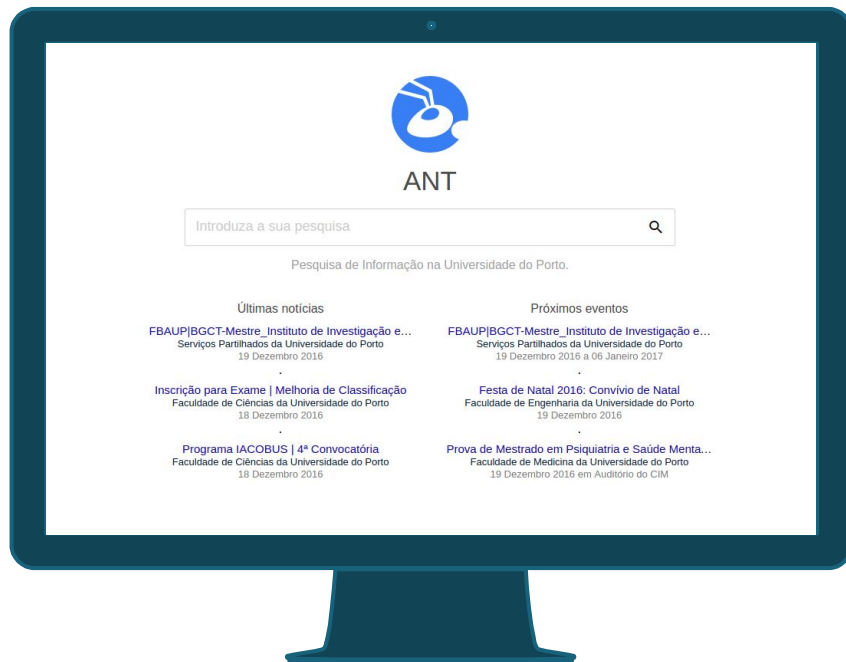
file pruning threshold.

Literature

Test collections

	Num. Docs.	Year(s)	Evaluation	Document ranking	Entity ranking	List completion
INEX 2009 Wikipedia Collection	2.6M XML	2009–2010	Offline	✓	✓	✓
TREC Washington Post Corpus	595K JSON	2018	Offline	✓	✗	✗
Social Science Open Access Repository	32K JSON	2017	Online	✓	✗	✗

Software





Contributions

- Graph-of-entity
- Hypergraph-of-entity

Example document

Semantic search seeks to improve search [Search Engine Technology] accuracy by understanding the searcher's intent [Intention] and the contextual [Contextual (language use)] meaning of terms as they appear in the searchable dataspace, whether on the Web [World Wide Web] or within a closed system, to generate more relevant results.

– ‘Semantic search’, Wikipedia, 9:10am, January 7, 2016

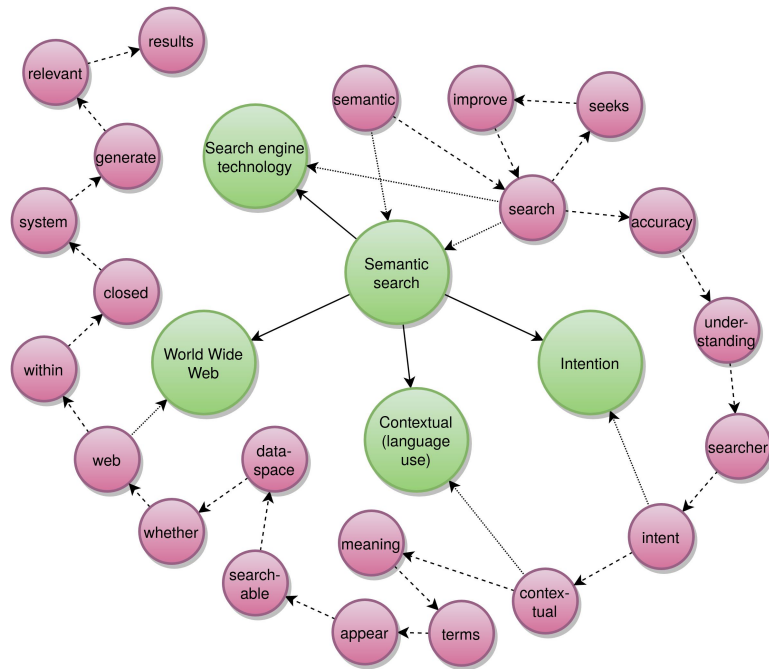
Graph-of-entity: representation

- Nodes

- 'term'
- 'entity'

- Edges (directed and unweighted)

- > 'before'
- 'related_to'
-> 'contained_in'



Note: The direction for 'related_to' edges has been corrected.

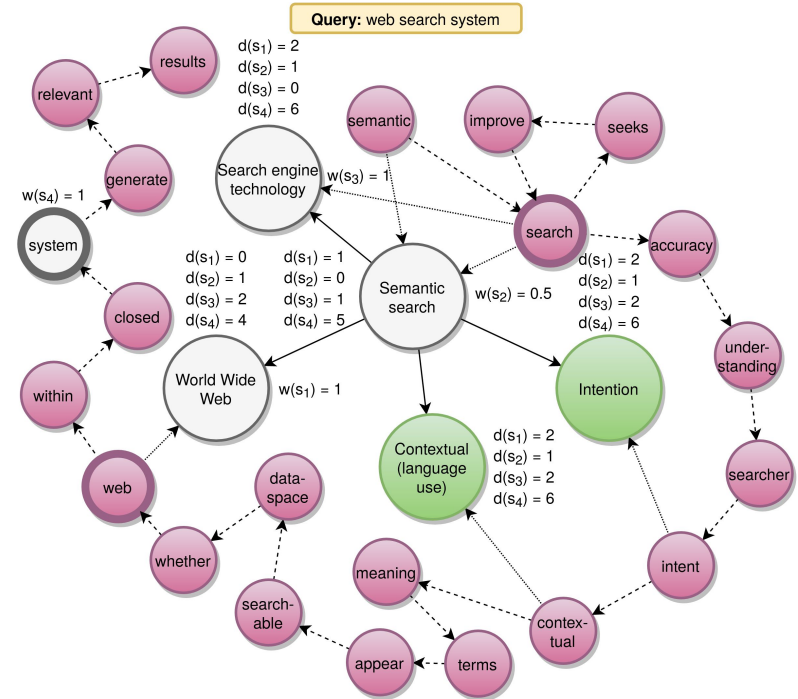
Seed nodes

- Map the query to the graph
- Can be expanded to adjacent entities
- Weighted according to their representability of the query

Graph-of-entity: retrieval

Three components:

- Coverage
- Confidence weight
- Entity weight



Note: The direction for 'related_to' edges has been corrected.

Scaling issues

INEX 2009 Wikipedia subset

- 7,484 documents
- Graph-of-entity
 - 981,647 nodes
 - 9,942,647 edges


Scaling issues

INEX 2009 Wikipedia subset

- 7,484 documents
- Graph-of-entity
 - 981,647 nodes
 - 9,942,647 edges


Scaling issues

INEX 2009 Wikipedia subset

- 7,484 documents
 - Graph-of-entity
 - 981,647 nodes
 - 9,942,647 edges
- 2 orders of magnitude
- 


Scaling issues

INEX 2009 Wikipedia subset

- 7,484 documents
 - Graph-of-entity
 - 981,647 nodes
 - 9,942,647 edges
- 3 orders of magnitude
- 

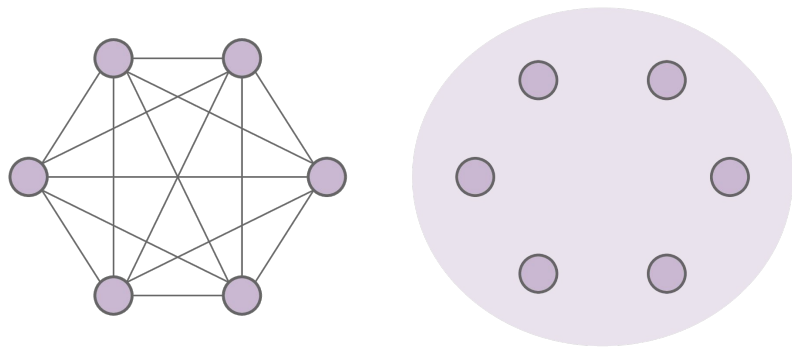
Scaling issues

INEX 2009 Wikipedia subset

- 7,484 documents
 - Graph-of-entity
 - 981,647 nodes
 - 9,942,647 edges
- 3 orders of magnitude
- 

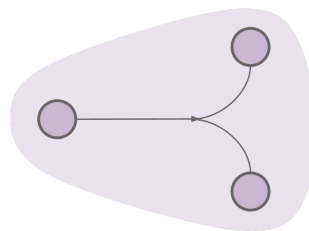
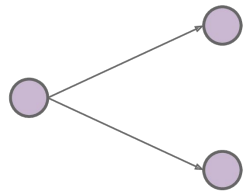
How could I reduce the number of edges per node?

From graphs to hypergraphs



Representing full connectivity

(e.g., synonyms)



Representing directed n-ary connectivity

(e.g., e-mail message)

Hypergraph-of-entity: joint representation model

Base model

■ Nodes

● 'term'

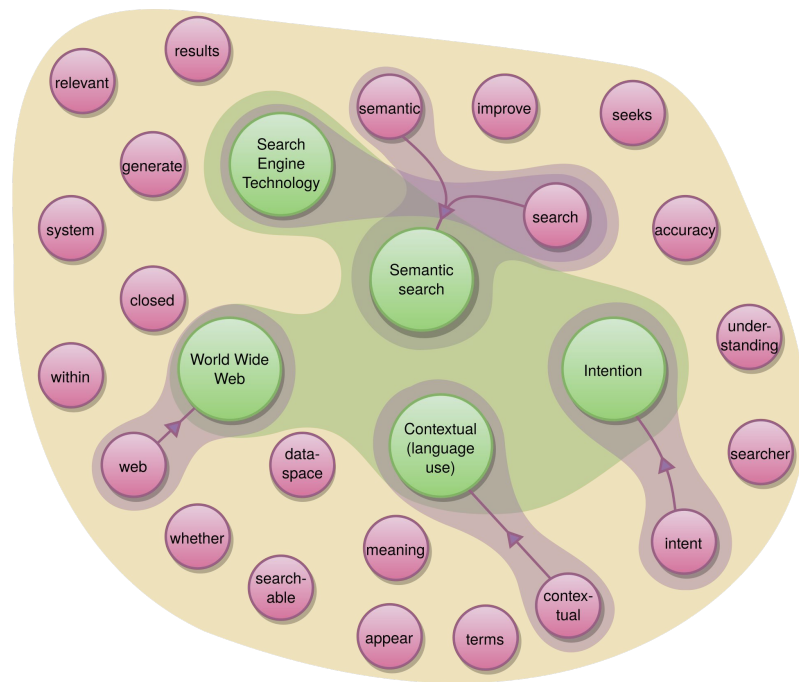
● 'entity'

■ Hyperedges

— 'document'

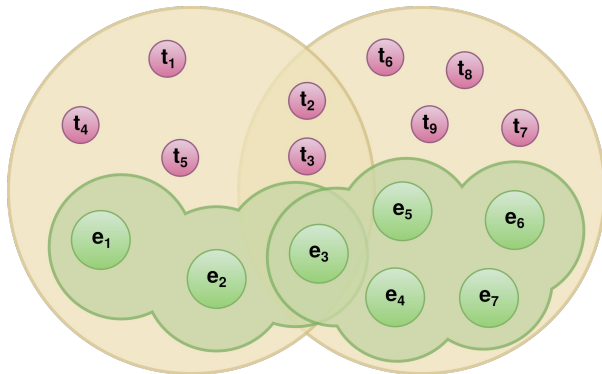
— 'related_to'

— 'contained_in'

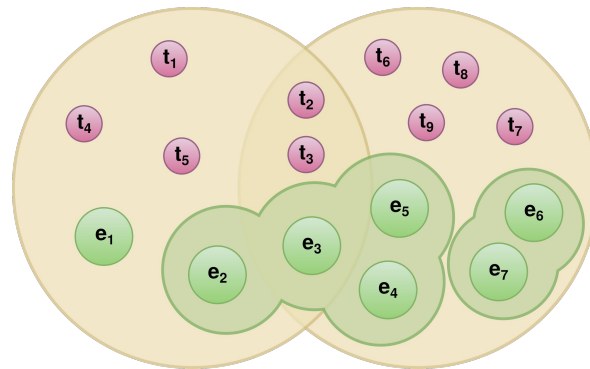


Hypergraph-of-entity: joint representation model

'related_to' hyperedges



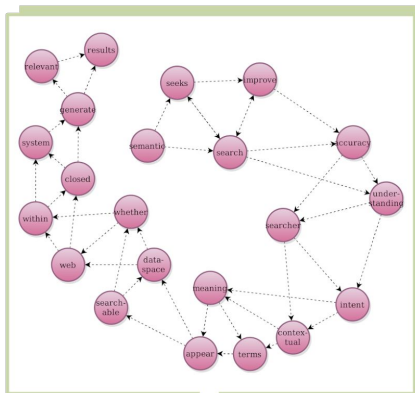
Grouped by co-occurrence



Grouped by subject

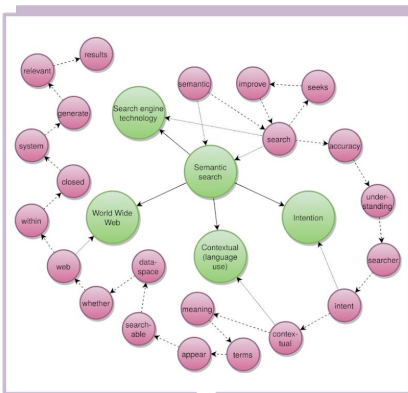
Scaling issues: mitigated

Graph-of-word



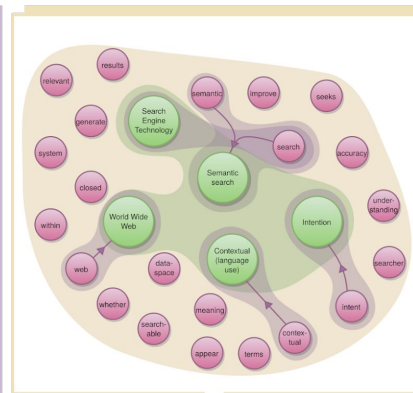
- 7,487 documents
- 492,185 vertices
- 22,906,803 edges
- $|E| = 46.5 \times |V|$

Graph-of-entity



- 7,487 documents
- 981,647 vertices
- 9,942,647 edges
- $|E| = 10.1 \times |V|$

Hypergraph-of-entity



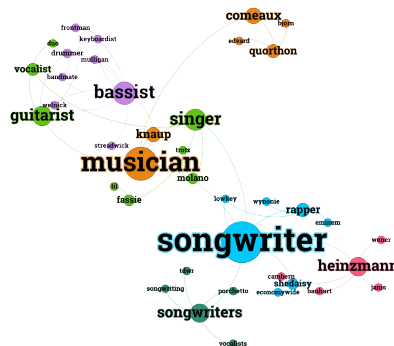
- 7,487 documents
- 607,213 vertices
- 253,154 hyperedges
- $|E| = 0.4 \times |V|$

Extensions

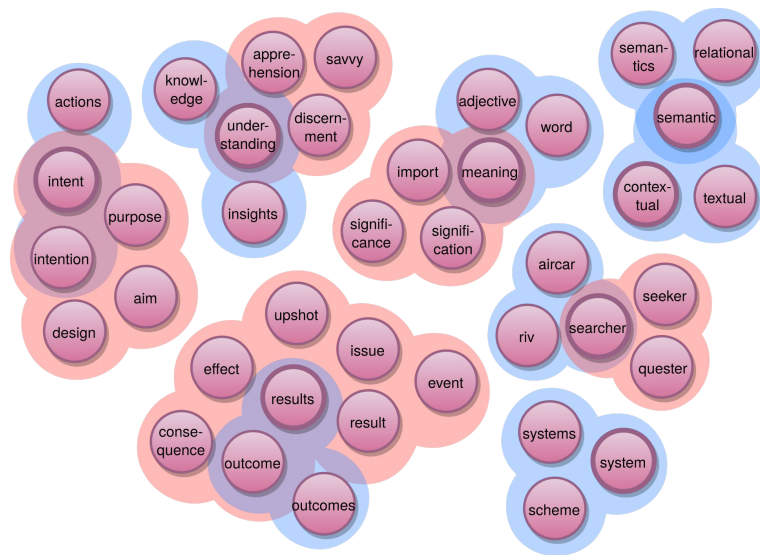
- WordNet 3.0 Nouns

- Word2vec SimNet

- size=100, window=5
- 2-NN, cosine similarity > 0.5
- Hyperedge per term neighborhood



Word2vec SimNet:
Ego Network for
'musician', with depth 3



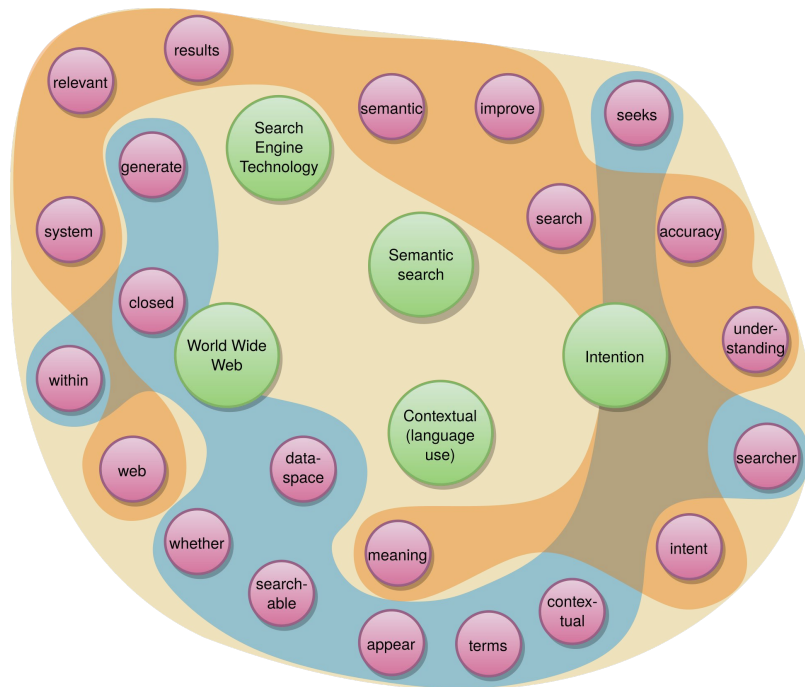
Synonyms + Context

Hypergraph-of-entity: joint representation model

Extensions

TF-bins (low TF  / high TF )

- Discretization of term frequency
- Optionally weighted by percentile order, e.g.:
 - For $P_2 = \{50, 100\}$
 - $w_{50} = \frac{1}{2}$, $w_{100} = 1$



TF-bins:
bin width = 2

Hypergraph-of-entity: joint representation model

Extensions

Weights

- To further constraint or guide the ranking function
- Add a bias that can affect both node and hyperedge sampling

Hypergraph-of-entity: universal ranking function

Random walk score

- Random walks on a hypergraph
- Launched from each seed node
- Final score computed as:
 - \sum weighted sum
 - Confidence weight
 - \times visitation probability
 - \times coverage

Hypergraph-of-entity: universal ranking function

Table 7.4: Mapping entity-oriented search tasks to the hypergraph-of-entity.

	Query	Input	Results	Output
Ad hoc document retrieval	Keyword	Term nodes	Documents	Hyperedge ranking
Ad hoc entity retrieval	Keyword	Term nodes	Entities	Node ranking
Related entity finding	Entity	One entity node	Entities	Node ranking
Entity list completion	Entity	Multiple entity nodes	Entities	Node ranking

Hypergraph-of-entity: universal ranking function

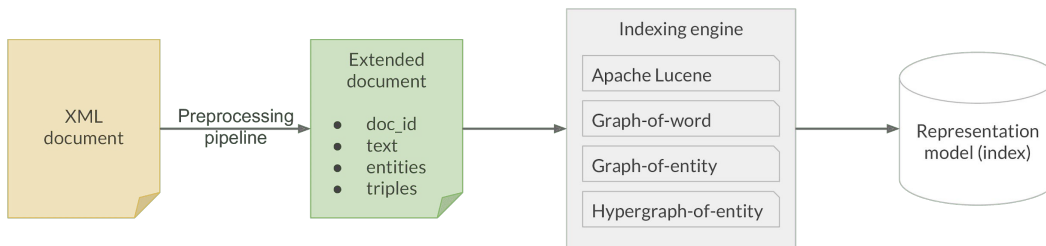
Table 7.5: Random walk score parameters and chosen configuration.

Parameter	Description	Configuration
ℓ	Length of the random walk.	2
r	Number of repeated random walks per seed node.	10,000
Δ_{nf}	Number of cycles of node fatigue (see Section B.2).	0
Δ_{ef}	Number of cycles of (hyper)edge fatigue (see Section B.2).	0
<i>expansion</i>	Whether to expand query to neighboring entities.	<i>false</i>
<i>directed</i>	Whether to consider or ignore direction.	<i>true</i>
<i>weighted</i>	Whether to consider node and hyperedge weights.	<i>false</i>

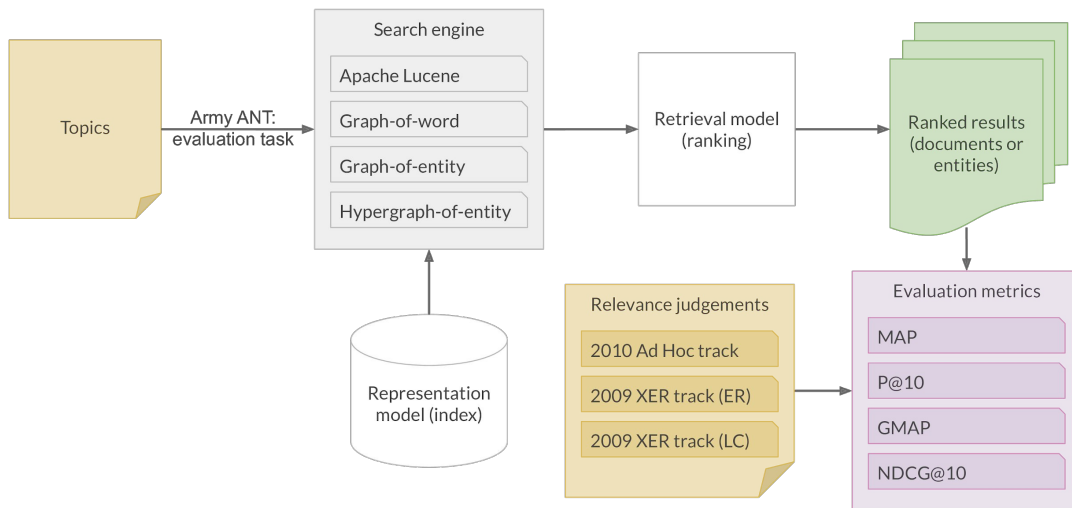
Evaluation: workflow

Indexing, retrieval and ranking, and evaluation

Indexing



Ranked retrieval & evaluation



Evaluation: main experiments

Retrieval performance over different representations:

- Text-only
- Base model
- Synonyms
- Context
- Syns+Cont.
- Cont.+Syns
- Syns+Cont.+Weights

Over different ranking function parameter configurations:

- Best results for:
 - Low ℓ
 - High r
 - No fatigue
- Variable results for:
 - Expansion
 - Weights

And for multiple tasks over the same index.

Evaluation: INEX 2009 10T-NL Wikipedia subset

Task

- Ad hoc document retrieval
2010 Ad Hoc track queries

Goal

- Compare graph-of-entity and
hypergraph-of-entity

Evaluation: INEX 2009 10T-NL Wikipedia subset

Ad hoc document retrieval Graph-of-entity vs hypergraph-of-entity

Table 9.5: Graph-of-entity (GoE) vs hypergraph-of-entity (HGoE) with $\ell = 2$.

(a) Effectiveness (highest values for Lucene and graph-based models in bold).

Index	Ranking	GMAP	MAP	Precision	Recall	NDCG@10	P@10
Lucene	TF-IDF	0.1540	0.1710	0.1389	0.8007	0.2671	0.2800
	BM25	0.2802	0.2963	0.1396	0.8241	0.5549	0.5000
GoE	EW	0.0003	0.0399	0.1771	0.2233	0.1480	0.1500
HGoE	RWS($r = 10^1$)	0.0000	0.0485	0.0734	0.3085	0.1229	0.1200
	RWS($r = 10^2$)	0.0546	0.1118	0.0342	0.7554	0.1474	0.1500
	RWS($r = 10^3$)	0.1017	0.1492	0.0199	0.9122	0.2074	0.2200
	RWS($r = 10^4$)	0.1224	0.1689	0.0167	0.9922	0.1699	0.1700

(b) Efficiency (lowest times for Lucene and graph-based models in bold).

Index	Ranking	Indexing Time (Total)	Search Time (Avg./Query)	Nodes	Edges
Lucene	TF-IDF	27s 769ms	209ms	N/A	N/A
	BM25		316ms		
GoE	EW	1h 38m	21s 557ms	981,647	9,942,647
HGoE	RWS($r = 10^1$)	53s 922ms	943ms	607,213	253,154
	RWS($r = 10^2$)		11s 134ms		
	RWS($r = 10^3$)		1m 17s 540ms		
	RWS($r = 10^4$)		13m 04s 057ms		

Evaluation: INEX 2009 10T-NL Wikipedia subset

Ad hoc document retrieval Graph-of-entity vs hypergraph-of-entity

Table 9.5: Graph-of-entity (GoE) vs hypergraph-of-entity (HGoE) with $\ell = 2$.

(a) Effectiveness (highest values for Lucene and graph-based models in bold).

Index	Ranking	GMAP	MAP	Precision	Recall	NDCG@10	P@10
Lucene	TF-IDF	0.1540	0.1710	0.1389	0.8007	0.2671	0.2800
	BM25	0.2802	0.2963	0.1396	0.8241	0.5549	0.5000
GoE	EW	0.0003	0.0399	0.1771	0.2233	0.1480	0.1500
HGoE	RWS($r = 10^1$)	0.0000	0.0485	0.0734	0.3085	0.1229	0.1200
	RWS($r = 10^2$)	0.0546	0.1118	0.0342	0.7554	0.1474	0.1500
	RWS($r = 10^3$)	0.1017	0.1492	0.0199	0.9122	0.2074	0.2200
	RWS($r = 10^4$)	0.1224	0.1689	0.0167	0.9922	0.1699	0.1700

(b) Efficiency (lowest times for Lucene and graph-based models in bold).

Index	Ranking	Indexing Time (Total)	Search Time (Avg./Query)	Nodes	Edges
Lucene	TF-IDF	27s 769ms	209ms	N/A	N/A
	BM25		316ms		
GoE	EW	1h 38m	21s 557ms	981,647	9,942,647
HGoE	RWS($r = 10^1$)	53s 922ms	943ms	607,213	253,154
	RWS($r = 10^2$)		11s 134ms		
	RWS($r = 10^3$)		1m 17s 540ms		
	RWS($r = 10^4$)		13m 04s 057ms		

Evaluation: INEX 2009 10T-NL Wikipedia subset

Ad hoc document retrieval Graph-of-entity vs hypergraph-of-entity

Table 9.5: Graph-of-entity (GoE) vs hypergraph-of-entity (HGoE) with $\ell = 2$.

(a) Effectiveness (highest values for Lucene and graph-based models in bold).

Index	Ranking	GMAP	MAP	Precision	Recall	NDCG@10	P@10
Lucene	TF-IDF	0.1540	0.1710	0.1389	0.8007	0.2671	0.2800
	BM25	0.2802	0.2963	0.1396	0.8241	0.5549	0.5000
GoE	EW	0.0003	0.0399	0.1771	0.2233	0.1480	0.1500
HGoE	RWS($r = 10^1$)	0.0000	0.0485	0.0734	0.3085	0.1229	0.1200
	RWS($r = 10^2$)	0.0546	0.1118	0.0342	0.7554	0.1474	0.1500
	RWS($r = 10^3$)	0.1017	0.1492	0.0199	0.9122	0.2074	0.2200
	RWS($r = 10^4$)	0.1224	0.1689	0.0167	0.9922	0.1699	0.1700

(b) Efficiency (lowest times for Lucene and graph-based models in bold).

Index	Ranking	Indexing Time (Total)	Search Time (Avg./Query)	Nodes	Edges
Lucene	TF-IDF	27s 769ms	209ms	N/A	N/A
	BM25		316ms		
GoE	EW	1h 38m	21s 557ms	981,647	9,942,647
HGoE	RWS($r = 10^1$)	53s 922ms	943ms	607,213	253,154
	RWS($r = 10^2$)		11s 134ms		
	RWS($r = 10^3$)		1m 17s 540ms		
	RWS($r = 10^4$)		13m 04s 057ms		

Evaluation: INEX 2009 10T-NL Wikipedia subset

Ad hoc document retrieval Graph-of-entity vs hypergraph-of-entity

Table 9.5: Graph-of-entity (GoE) vs hypergraph-of-entity (HGoE) with $\ell = 2$.

(a) Effectiveness (highest values for Lucene and graph-based models in bold).

Index	Ranking	GMAP	MAP	Precision	Recall	NDCG@10	P@10
Lucene	TF-IDF	0.1540	0.1710	0.1389	0.8007	0.2671	0.2800
	BM25	0.2802	0.2963	0.1396	0.8241	0.5549	0.5000
GoE	EW	0.0003	0.0399	0.1771	0.2233	0.1480	0.1500
HGoE	RWS($r = 10^1$)	0.0000	0.0485	0.0734	0.3085	0.1229	0.1200
	RWS($r = 10^2$)	0.0546	0.1118	0.0342	0.7554	0.1474	0.1500
	RWS($r = 10^3$)	0.1017	0.1492	0.0199	0.9122	0.2074	0.2200
	RWS($r = 10^4$)	0.1224	0.1689	0.0167	0.9922	0.1699	0.1700

(b) Efficiency (lowest times for Lucene and graph-based models in bold).

Index	Ranking	Indexing Time (Total)	Search Time (Avg./Query)	Nodes	Edges
Lucene	TF-IDF	27s 769ms	209ms	N/A	N/A
	BM25		316ms		
GoE	EW	1h 38m	21s 557ms	981,647	9,942,647
HGoE	RWS($r = 10^1$)	53s 922ms	943ms	607,213	253,154
	RWS($r = 10^2$)		11s 134ms		
	RWS($r = 10^3$)		1m 17s 540ms		
	RWS($r = 10^4$)		13m 04s 057ms		

Evaluation: INEX 2009 52T-NL Wikipedia subset

Task

- Ad hoc document retrieval
2010 Ad Hoc track qrels

Goal

- Compare representation models
based on the hypergraph-of-entity

Evaluation: INEX 2009 52T-NL Wikipedia subset

Ad hoc document retrieval Comparing representation models

Table 9.4: Best overall parameter configuration according to the mean average precision.

(a) Effectiveness (highest values for Lucene and hypergraph-of-entity in bold; differences in MAP are not statistically significant, except between the Lucene baselines and the hypergraph-of-entity indexes).

Index	Ranking	GMAP	MAP	Precision	Recall	NDCG@10	P@10
Lucene	TF-IDF	0.1345	0.1689	0.0650	0.8476	0.2291	0.2346
	BM25	0.2740	0.3269	0.0647	0.8598	0.5607	0.5250
Hypergraph-of-Entity: Random Walk Score ($\ell = 2$, $r = 10^3$)							
<i>Base Model</i>	RWS	0.0285	0.0864	0.0219	0.8003	0.1413	0.1269
<i>Syns</i>	RWS	0.0281	0.0840	0.0225	0.8099	0.1301	0.1231
<i>Context</i>	RWS	0.0134	0.0811	0.0220	0.8027	0.1218	0.1192
<i>Syns+Context</i>	RWS	0.0299	0.0837	0.0236	0.8069	0.1310	0.1231
<i>Context+Syns</i>	RWS	0.0296	0.0814	0.0242	0.8148	0.1256	0.1250
<i>Syns+Cont.+Weights</i>	RWS	0.0313	0.0884	0.0274	0.8059	0.1256	0.1154

(b) Efficiency (lowest times for Lucene and hypergraph-of-entity in bold).

Index	Ranking	Indexing Time		Search Time	
		Avg./Doc	Total	Avg./Query	Total
Lucene	TF-IDF	2.16ms	1m 21s 382ms	1s 148ms	59s 698ms
	BM25			1s 220ms	1m 03s 461ms
Hypergraph-of-Entity: Random Walk Score ($\ell = 2$, $r = 10^3$)					
<i>Base Model</i>	RWS	6.52ms	4m 05s 612ms	3m 22s 826ms	2h 55m 47s
<i>Syns</i>	RWS	6.22ms	3m 54s 587ms	3m 31s 038ms	3h 02m 54s
<i>Context</i>	RWS	6.35ms	3m 59s 446ms	3m 35s 623ms	3h 06m 52s
<i>Syns+Context</i>	RWS	6.29ms	3m 57s 264ms	3m 33s 000ms	3h 04m 36s
<i>Context+Syns</i>	RWS	6.33ms	3m 58s 659ms	3m 36s 487ms	3h 07m 37s
<i>Syns+Cont.+Weights</i>	RWS	6.52ms	4m 05s 984ms	10m 55s 590ms	9h 28m 11s

Evaluation: INEX 2009 52T-NL Wikipedia subset

Ad hoc document retrieval Comparing representation models

Table 9.4: Best overall parameter configuration according to the mean average precision.

(a) Effectiveness (highest values for Lucene and hypergraph-of-entity in bold; differences in MAP are not statistically significant, except between the Lucene baselines and the hypergraph-of-entity indexes).

Index	Ranking	GMAP	MAP	Precision	Recall	NDCG@10	P@10
Lucene	TF-IDF	0.1345	0.1689	0.0650	0.8476	0.2291	0.2346
	BM25	0.2740	0.3269	0.0647	0.8598	0.5607	0.5250
Hypergraph-of-Entity: Random Walk Score ($\ell = 2$, $r = 10^3$)							
<i>Base Model</i>	RWS	0.0285	0.0864	0.0219	0.8003	0.1413	0.1269
<i>Syns</i>	RWS	0.0281	0.0840	0.0225	0.8099	0.1301	0.1231
<i>Context</i>	RWS	0.0134	0.0811	0.0220	0.8027	0.1218	0.1192
<i>Syns+Context</i>	RWS	0.0299	0.0837	0.0236	0.8069	0.1310	0.1231
<i>Context+Syns</i>	RWS	0.0296	0.0814	0.0242	0.8148	0.1256	0.1250
<i>Syns+Cont.+Weights</i>	RWS	0.0313	0.0884	0.0274	0.8059	0.1256	0.1154

(b) Efficiency (lowest times for Lucene and hypergraph-of-entity in bold).

Index	Ranking	Indexing Time		Search Time	
		Avg./Doc	Total	Avg./Query	Total
Lucene	TF-IDF	2.16ms	1m 21s 382ms	1s 148ms	59s 698ms
	BM25			1s 220ms	1m 03s 461ms
Hypergraph-of-Entity: Random Walk Score ($\ell = 2$, $r = 10^3$)					
<i>Base Model</i>	RWS	6.52ms	4m 05s 612ms	3m 22s 826ms	2h 55m 47s
<i>Syns</i>	RWS	6.22ms	3m 54s 587ms	3m 31s 038ms	3h 02m 54s
<i>Context</i>	RWS	6.35ms	3m 59s 446ms	3m 35s 623ms	3h 06m 52s
<i>Syns+Context</i>	RWS	6.29ms	3m 57s 264ms	3m 33s 000ms	3h 04m 36s
<i>Context+Syns</i>	RWS	6.33ms	3m 58s 659ms	3m 36s 487ms	3h 07m 37s
<i>Syns+Cont.+Weights</i>	RWS	6.52ms	4m 05s 984ms	10m 55s 590ms	9h 28m 11s

Evaluation: INEX 2009 52T-NL Wikipedia subset

Ad hoc document retrieval Comparing representation models

Table 9.4: Best overall parameter configuration according to the mean average precision.

(a) Effectiveness (highest values for Lucene and hypergraph-of-entity in bold; differences in MAP are not statistically significant, except between the Lucene baselines and the hypergraph-of-entity indexes).

Index	Ranking	GMAP	MAP	Precision	Recall	NDCG@10	P@10
Lucene	TF-IDF	0.1345	0.1689	0.0650	0.8476	0.2291	0.2346
	BM25	0.2740	0.3269	0.0647	0.8598	0.5607	0.5250
Hypergraph-of-Entity: Random Walk Score ($\ell = 2$, $r = 10^3$)							
<i>Base Model</i>	RWS	0.0285	0.0864	0.0219	0.8003	0.1413	0.1269
<i>Syns</i>	RWS	0.0281	0.0840	0.0225	0.8099	0.1301	0.1231
<i>Context</i>	RWS	0.0134	0.0811	0.0220	0.8027	0.1218	0.1192
<i>Syns+Context</i>	RWS	0.0299	0.0837	0.0236	0.8069	0.1310	0.1231
<i>Context+Syns</i>	RWS	0.0296	0.0814	0.0242	0.8148	0.1256	0.1250
<i>Syns+Cont.+Weights</i>	RWS	0.0313	0.0884	0.0274	0.8059	0.1256	0.1154

(b) Efficiency (lowest times for Lucene and hypergraph-of-entity in bold).

Index	Ranking	Indexing Time		Search Time	
		Avg./Doc	Total	Avg./Query	Total
Lucene	TF-IDF	2.16ms	1m 21s 382ms	1s 148ms	59s 698ms
	BM25			1s 220ms	1m 03s 461ms
Hypergraph-of-Entity: Random Walk Score ($\ell = 2$, $r = 10^3$)					
Base Model	RWS	6.52ms	4m 05s 612ms	3m 22s 826ms	2h 55m 47s
Syns	RWS	6.22ms	3m 54s 587ms	3m 31s 038ms	3h 02m 54s
Context	RWS	6.35ms	3m 59s 446ms	3m 35s 623ms	3h 06m 52s
Syns+Context	RWS	6.29ms	3m 57s 264ms	3m 33s 000ms	3h 04m 36s
Context+Syns	RWS	6.33ms	3m 58s 659ms	3m 36s 487ms	3h 07m 37s
Syns+Contl.+Weights	RWS	6.52ms	4m 05s 984ms	10m 55s 590ms	9h 28m 11s

Evaluation: INEX 2009 52T-NL Wikipedia subset

Ad hoc document retrieval Comparing representation models

Table 9.4: Best overall parameter configuration according to the mean average precision.

(a) Effectiveness (highest values for Lucene and hypergraph-of-entity in bold; differences in MAP are not statistically significant, except between the Lucene baselines and the hypergraph-of-entity indexes).

Index	Ranking	GMAP	MAP	Precision	Recall	NDCG@10	P@10
Lucene	TF-IDF	0.1345	0.1689	0.0650	0.8476	0.2291	0.2346
	BM25	0.2740	0.3269	0.0647	0.8598	0.5607	0.5250
Hypergraph-of-Entity: Random Walk Score ($\ell = 2, r = 10^3$)							
<i>Base Model</i>	RWS	0.0285	0.0864	0.0219	0.8003	0.1413	0.1269
<i>Syns</i>	RWS	0.0281	0.0840	0.0225	0.8099	0.1301	0.1231
<i>Context</i>	RWS	0.0134	0.0811	0.0220	0.8027	0.1218	0.1192
<i>Syns+Context</i>	RWS	0.0299	0.0837	0.0236	0.8069	0.1310	0.1231
<i>Context+Syns</i>	RWS	0.0296	0.0814	0.0242	0.8148	0.1256	0.1250
<i>Syns+Cont.+Weights</i>	RWS	0.0313	0.0884	0.0274	0.8059	0.1256	0.1154

(b) Efficiency (lowest times for Lucene and hypergraph-of-entity in bold).

Index	Ranking	Indexing Time		Search Time	
		Avg./Doc	Total	Avg./Query	Total
Lucene	TF-IDF	2.16ms	1m 21s 382ms	1s 148ms	59s 698ms
	BM25			1s 220ms	1m 03s 461ms
Hypergraph-of-Entity: Random Walk Score ($\ell = 2, r = 10^3$)					
Base Model	RWS	6.52ms	4m 05s 612ms	3m 22s 826ms	2h 55m 47s
Syns	RWS	6.22ms	3m 54s 587ms	3m 31s 038ms	3h 02m 54s
Context	RWS	6.35ms	3m 59s 446ms	3m 35s 623ms	3h 06m 52s
Syns+Context	RWS	6.29ms	3m 57s 264ms	3m 33s 000ms	3h 04m 36s
Context+Syns	RWS	6.33ms	3m 58s 659ms	3m 36s 487ms	3h 07m 37s
Syns+Cont.+Weights	RWS	6.52ms	4m 05s 984ms	10m 55s 590ms	9h 28m 11s

Evaluation: INEX 2009 Wikipedia full collection

Tasks

- Ad hoc document retrieval
2010 Ad Hoc track qrels
- Ad hoc entity retrieval
2009 XER track entity ranking qrels
- Entity list completion
2009 XER track list completion qrels

Goal

- Assess the viability of a general retrieval model

Reduced representation

- Keyword-based document profiles

Evaluation: INEX 2009 Wikipedia full collection

Multiple tasks

Keyword-based document profiles

Note: The entity index was created from sentence-based entity profiles

Table 10.6: Overall comparison of retrieval performance, for multiple entity-oriented search tasks, based on the complete INEX 2009 Wikipedia collection.

Index	Ranking	MAP	P@10	Index Time	Avg./Query	Nodes	Edges
AD HOC DOCUMENT RETRIEVAL							
Lucene	TF-IDF	0.0228	0.0692	15h 06m	460ms	—	—
	BM25	0.0324	0.1173		370ms	—	—
Fixed parameters over HGoE variations: $RWS(\ell = 2, r = 10^4, \Delta_{nf} = 0, \Delta_{ef} = 0, exp = F, wei = F)$							
Base Model	RWS	0.0863	0.2462	33h 53m	23s 405ms	3,506,823	7,721,743
Syns	RWS	0.0937	0.2615	33h 05m	55s 555ms	3,510,462	7,734,931
Cont.	RWS	0.0869	0.2654	34h 37m	24s 348ms	3,604,185	7,929,841
TF-Bins ₂	RWS	0.0172	0.0500	35h 26m	2m 58s	3,506,823	10,338,867
Syns+Cont.	RWS	0.0882	0.2692	37h 16m	23s 265ms	3,606,693	7,945,083
AD HOC ENTITY RETRIEVAL							
Lucene	TF-IDF	0.0373	0.0636	59h 17m	1s 370ms	—	—
	BM25	0.0668	0.1182		798ms	—	—
Fixed parameters over HGoE variations: $RWS(\ell = 2, r = 10^4, \Delta_{nf} = 0, \Delta_{ef} = 0, exp = F, wei = F)$							
Base Model	RWS	0.1390	0.2455	33h 53m	26s 330ms	3,506,823	7,721,743
Syns	RWS	0.1337	0.2473	33h 05m	30s 232ms	3,510,462	7,734,931
Cont.	RWS	0.1304	0.2364	34h 37m	27s 620ms	3,604,185	7,929,841
TF-Bins ₂	RWS	0.0300	0.1145	35h 26m	4m 41s	3,506,823	10,338,867
Syns+Cont.	RWS	0.1313	0.2509	37h 16m	26s 877ms	3,606,693	7,945,083
ENTITY LIST COMPLETION							
Lucene	TF-IDF	0.0558	0.1000	59h 17m	1s 230ms	—	—
	BM25	0.0666	0.1250		1s 221ms	—	—
Fixed parameters over HGoE variations: $RWS(\ell = 2, r = 10^4, \Delta_{nf} = 0, \Delta_{ef} = 0, exp = F, wei = F)$							
Base Model	RWS	0.0879	0.0769	33h 53m	19s 162ms	3,506,823	7,721,743
Syns	RWS	0.0857	0.0635	33h 05m	19s 875ms	3,510,462	7,734,931
Cont.	RWS	0.0875	0.0692	34h 37m	19s 422ms	3,604,185	7,929,841
TF-Bins ₂	RWS	0.0006	0.0058	35h 26m	1m 08s	3,506,823	10,338,867
Syns+Cont.	RWS	0.0884	0.0788	37h 16m	19s 824ms	3,606,693	7,945,083

Evaluation: INEX 2009 Wikipedia full collection

Multiple tasks

Keyword-based document profiles

Note: The entity index was created from sentence-based entity profiles

Table 10.6: Overall comparison of retrieval performance, for multiple entity-oriented search tasks, based on the complete INEX 2009 Wikipedia collection.

Index	Ranking	MAP	P@10	Index Time	Avg./Query	Nodes	Edges
AD HOC DOCUMENT RETRIEVAL							
Lucene	TF-IDF	0.0228	0.0692	15h 06m	460ms	—	—
	BM25	0.0324	0.1173		370ms	—	—
Fixed parameters over HGoE variations: $RWS(\ell = 2, r = 10^4, \Delta_{nf} = 0, \Delta_{ef} = 0, exp = F, wei = F)$							
Base Model	RWS	0.0863	0.2462	33h 53m	23s 405ms	3,506,823	7,721,743
Syns	RWS	0.0937	0.2615	33h 05m	55s 555ms	3,510,462	7,734,931
Cont.	RWS	0.0869	0.2654	34h 37m	24s 348ms	3,604,185	7,929,841
TF-Bins ₂	RWS	0.0172	0.0500	35h 26m	2m 58s	3,506,823	10,338,867
Syns+Cont.	RWS	0.0882	0.2692	37h 16m	23s 265ms	3,606,693	7,945,083
AD HOC ENTITY RETRIEVAL							
Lucene	TF-IDF	0.0373	0.0636	59h 17m	1s 370ms	—	—
	BM25	0.0668	0.1182		798ms	—	—
Fixed parameters over HGoE variations: $RWS(\ell = 2, r = 10^4, \Delta_{nf} = 0, \Delta_{ef} = 0, exp = F, wei = F)$							
Base Model	RWS	0.1390	0.2455	33h 53m	26s 330ms	3,506,823	7,721,743
Syns	RWS	0.1337	0.2473	33h 05m	30s 232ms	3,510,462	7,734,931
Cont.	RWS	0.1304	0.2364	34h 37m	27s 620ms	3,604,185	7,929,841
TF-Bins ₂	RWS	0.0300	0.1145	35h 26m	4m 41s	3,506,823	10,338,867
Syns+Cont.	RWS	0.1313	0.2509	37h 16m	26s 877ms	3,606,693	7,945,083
ENTITY LIST COMPLETION							
Lucene	TF-IDF	0.0558	0.1000	59h 17m	1s 230ms	—	—
	BM25	0.0666	0.1250		1s 221ms	—	—
Fixed parameters over HGoE variations: $RWS(\ell = 2, r = 10^4, \Delta_{nf} = 0, \Delta_{ef} = 0, exp = F, wei = F)$							
Base Model	RWS	0.0879	0.0769	33h 53m	19s 162ms	3,506,823	7,721,743
Syns	RWS	0.0857	0.0635	33h 05m	19s 875ms	3,510,462	7,734,931
Cont.	RWS	0.0875	0.0692	34h 37m	19s 422ms	3,604,185	7,929,841
TF-Bins ₂	RWS	0.0006	0.0058	35h 26m	1m 08s	3,506,823	10,338,867
Syns+Cont.	RWS	0.0884	0.0788	37h 16m	19s 824ms	3,606,693	7,945,083

Evaluation: INEX 2009 Wikipedia full collection

Multiple tasks

Keyword-based document profiles

Note: The entity index was created from sentence-based entity profiles

Table 10.6: Overall comparison of retrieval performance, for multiple entity-oriented search tasks, based on the complete INEX 2009 Wikipedia collection.

Index	Ranking	MAP	P@10	Index Time	Avg./Query	Nodes	Edges
AD HOC DOCUMENT RETRIEVAL							
Lucene	TF-IDF	0.0228	0.0692	15h 06m	460ms	—	—
	BM25	0.0324	0.1173		370ms	—	—
Fixed parameters over HGoE variations: $RWS(\ell = 2, r = 10^4, \Delta_{nf} = 0, \Delta_{ef} = 0, exp = F, wei = F)$							
Base Model	RWS	0.0863	0.2462	33h 53m	23s 405ms	3,506,823	7,721,743
Syns	RWS	0.0937	0.2615	33h 05m	55s 555ms	3,510,462	7,734,931
Cont.	RWS	0.0869	0.2654	34h 37m	24s 348ms	3,604,185	7,929,841
TF-Bins ₂	RWS	0.0172	0.0500	35h 26m	2m 58s	3,506,823	10,338,867
Syns+Cont.	RWS	0.0882	0.2692	37h 16m	23s 265ms	3,606,693	7,945,083
AD HOC ENTITY RETRIEVAL							
Lucene	TF-IDF	0.0373	0.0636	59h 17m	1s 370ms	—	—
	BM25	0.0668	0.1182		798ms	—	—
Fixed parameters over HGoE variations: $RWS(\ell = 2, r = 10^4, \Delta_{nf} = 0, \Delta_{ef} = 0, exp = F, wei = F)$							
Base Model	RWS	0.1390	0.2455	33h 53m	26s 330ms	3,506,823	7,721,743
Syns	RWS	0.1337	0.2473	33h 05m	30s 232ms	3,510,462	7,734,931
Cont.	RWS	0.1304	0.2364	34h 37m	27s 620ms	3,604,185	7,929,841
TF-Bins ₂	RWS	0.0300	0.1145	35h 26m	4m 41s	3,506,823	10,338,867
Syns+Cont.	RWS	0.1313	0.2509	37h 16m	26s 877ms	3,606,693	7,945,083
ENTITY LIST COMPLETION							
Lucene	TF-IDF	0.0558	0.1000	59h 17m	1s 230ms	—	—
	BM25	0.0666	0.1250		1s 221ms	—	—
Fixed parameters over HGoE variations: $RWS(\ell = 2, r = 10^4, \Delta_{nf} = 0, \Delta_{ef} = 0, exp = F, wei = F)$							
Base Model	RWS	0.0879	0.0769	33h 53m	19s 162ms	3,506,823	7,721,743
Syns	RWS	0.0857	0.0635	33h 05m	19s 875ms	3,510,462	7,734,931
Cont.	RWS	0.0875	0.0692	34h 37m	19s 422ms	3,604,185	7,929,841
TF-Bins ₂	RWS	0.0006	0.0058	35h 26m	1m 08s	3,506,823	10,338,867
Syns+Cont.	RWS	0.0884	0.0788	37h 16m	19s 824ms	3,606,693	7,945,083

Evaluation: INEX 2009 Wikipedia full collection

Multiple tasks

Keyword-based document profiles

Note: The entity index was created from sentence-based entity profiles

Table 10.6: Overall comparison of retrieval performance, for multiple entity-oriented search tasks, based on the complete INEX 2009 Wikipedia collection.

Index	Ranking	MAP	P@10	Index Time	Avg./Query	Nodes	Edges
AD HOC DOCUMENT RETRIEVAL							
Lucene	TF-IDF	0.0228	0.0692	15h 06m	460ms	—	—
	BM25	0.0324	0.1173		370ms	—	—
Fixed parameters over HGoE variations: $RWS(\ell = 2, r = 10^4, \Delta_{nf} = 0, \Delta_{ef} = 0, exp = F, wei = F)$							
Base Model	RWS	0.0863	0.2462	33h 53m	23s 405ms	3,506,823	7,721,743
Syns	RWS	0.0937	0.2615	33h 05m	55s 555ms	3,510,462	7,734,931
Cont.	RWS	0.0869	0.2654	34h 37m	24s 348ms	3,604,185	7,929,841
TF-Bins ₂	RWS	0.0172	0.0500	35h 26m	2m 58s	3,506,823	10,338,867
Syns+Cont.	RWS	0.0882	0.2692	37h 16m	23s 265ms	3,606,693	7,945,083
AD HOC ENTITY RETRIEVAL							
Lucene	TF-IDF	0.0373	0.0636	59h 17m	1s 370ms	—	—
	BM25	0.0668	0.1182		798ms	—	—
Fixed parameters over HGoE variations: $RWS(\ell = 2, r = 10^4, \Delta_{nf} = 0, \Delta_{ef} = 0, exp = F, wei = F)$							
Base Model	RWS	0.1390	0.2455	33h 53m	26s 330ms	3,506,823	7,721,743
Syns	RWS	0.1337	0.2473	33h 05m	30s 232ms	3,510,462	7,734,931
Cont.	RWS	0.1304	0.2364	34h 37m	27s 620ms	3,604,185	7,929,841
TF-Bins ₂	RWS	0.0300	0.1145	35h 26m	4m 41s	3,506,823	10,338,867
Syns+Cont.	RWS	0.1313	0.2509	37h 16m	26s 877ms	3,606,693	7,945,083
ENTITY LIST COMPLETION							
Lucene	TF-IDF	0.0558	0.1000	59h 17m	1s 230ms	—	—
	BM25	0.0666	0.1250		1s 221ms	—	—
Fixed parameters over HGoE variations: $RWS(\ell = 2, r = 10^4, \Delta_{nf} = 0, \Delta_{ef} = 0, exp = F, wei = F)$							
Base Model	RWS	0.0879	0.0769	33h 53m	19s 162ms	3,506,823	7,721,743
Syns	RWS	0.0857	0.0635	33h 05m	19s 875ms	3,510,462	7,734,931
Cont.	RWS	0.0875	0.0692	34h 37m	19s 422ms	3,604,185	7,929,841
TF-Bins ₂	RWS	0.0006	0.0058	35h 26m	1m 08s	3,506,823	10,338,867
Syns+Cont.	RWS	0.0884	0.0788	37h 16m	19s 824ms	3,606,693	7,945,083

Evaluation: INEX 2009 Wikipedia full collection

Multiple tasks

Keyword-based document profiles

Note: The entity index was created from sentence-based entity profiles

Table 10.6: Overall comparison of retrieval performance, for multiple entity-oriented search tasks, based on the complete INEX 2009 Wikipedia collection.

Index	Ranking	MAP	P@10	Index Time	Avg./Query	Nodes	Edges
AD HOC DOCUMENT RETRIEVAL							
Lucene	TF-IDF	0.0228	0.0692	15h 06m	460ms	—	—
	BM25	0.0324	0.1173		370ms	—	—
Fixed parameters over HGoE variations: $RWS(\ell = 2, r = 10^4, \Delta_{nf} = 0, \Delta_{ef} = 0, exp = F, wei = F)$							
Base Model	RWS	0.0863	0.2462	33h 53m	23s 405ms	3,506,823	7,721,743
Syns	RWS	0.0937	0.2615	33h 05m	55s 555ms	3,510,462	7,734,931
Cont.	RWS	0.0869	0.2654	34h 37m	24s 348ms	3,604,185	7,929,841
TF-Bins ₂	RWS	0.0172	0.0500	35h 26m	2m 58s	3,506,823	10,338,867
Syns+Cont.	RWS	0.0882	0.2692	37h 16m	23s 265ms	3,606,693	7,945,083
AD HOC ENTITY RETRIEVAL							
Lucene	TF-IDF	0.0373	0.0636	59h 17m	1s 370ms	—	—
	BM25	0.0668	0.1182		798ms	—	—
Fixed parameters over HGoE variations: $RWS(\ell = 2, r = 10^4, \Delta_{nf} = 0, \Delta_{ef} = 0, exp = F, wei = F)$							
Base Model	RWS	0.1390	0.2455	33h 53m	26s 330ms	3,506,823	7,721,743
Syns	RWS	0.1337	0.2473	33h 05m	30s 232ms	3,510,462	7,734,931
Cont.	RWS	0.1304	0.2364	34h 37m	27s 620ms	3,604,185	7,929,841
TF-Bins ₂	RWS	0.0300	0.1145	35h 26m	4m 41s	3,506,823	10,338,867
Syns+Cont.	RWS	0.1313	0.2509	37h 16m	26s 877ms	3,606,693	7,945,083
ENTITY LIST COMPLETION							
Lucene	TF-IDF	0.0558	0.1000	59h 17m	1s 230ms	—	—
	BM25	0.0666	0.1250		1s 221ms	—	—
Fixed parameters over HGoE variations: $RWS(\ell = 2, r = 10^4, \Delta_{nf} = 0, \Delta_{ef} = 0, exp = F, wei = F)$							
Base Model	RWS	0.0879	0.0769	33h 53m	19s 162ms	3,506,823	7,721,743
Syns	RWS	0.0857	0.0635	33h 05m	19s 875ms	3,510,462	7,734,931
Cont.	RWS	0.0875	0.0692	34h 37m	19s 422ms	3,604,185	7,929,841
TF-Bins ₂	RWS	0.0006	0.0058	35h 26m	1m 08s	3,506,823	10,338,867
Syns+Cont.	RWS	0.0884	0.0788	37h 16m	19s 824ms	3,606,693	7,945,083

IV

Conclusions

Discussion, final remarks and future work

Discussion

- Efficiency / effectiveness trade-off:
 - Lower r is more efficient
 - Higher r is more effective
- Current implementation is:
 - Less efficient overall when compared to Lucene
 - More effective in the experiments using keyword-based document profiles

Discussion

- The hypergraph-of-entity is a novel indexing model for EOS
- The random walk score is the first attempt at a universal ranking function
 - Able to approximate or even beat TF-IDF and BM25
- Performance was better for:
 - Small datasets
 - Reduced document representations

Final remarks

- I have proven that a graph-based model is viable in EOS...
 - ...as a joint representation of corpora and knowledge bases...
 - ...using a universal ranking function to solve multiple EOS tasks.
-
- I improved retrieval effectiveness in some particular cases...
 - ...motivating the continued research of hypergraph-based models...
 - ...and unified frameworks in information retrieval.

Future work

- Optimize overall performance:
 - Prune nodes and hyperedges
 - Test different weighting functions
- Explore algebraic approaches:
 - Tensor-based representation
 - Personalized multilinear PageRank
- Expand supported tasks, e.g.:
 - Personalized search ('user' hyperedges)

THANK YOU!

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

<https://github.com/feup-infolab/army-ant/>

- J** J. Devezas and S. Nunes. "Characterizing the hypergraph-of-entity and the structural impact of its extensions". In: Applied Network Science - Special Issue of the 8th International Conference on Complex Networks and Their Applications 5:1 (2020), p. 79. issn: 2364-8228. doi: 10.1007/s41109-020-00320-z
- C** J. Devezas. "Graph-Based Entity-Oriented Search: A Unified Framework in Information Retrieval". In: Advances in Information Retrieval - 42nd European Conference on IR Research, ECIR 2020, Lisbon, Portugal, April 14-17, 2020, Proceedings, Part II, ed. by J. M. Jose, E. Yilmaz, J. Magalhães, P. Castells, N. Ferro, M. J. Silva, and F. Martins. Vol. 12036. Lecture Notes in Computer Science. Springer, 2020, pp. 602-607. doi: 10.1007/978-3-030-45442-5_78
- C** J. L. Devezas and S. Nunes. "Army ANT: A Workbench for Innovation in Entity-Oriented Search". In: Advances in Information Retrieval - 42nd European Conference on IR Research, ECIR 2020, Lisbon, Portugal, April 14-17, 2020, Proceedings, Part II, ed. by J. M. Jose, E. Yilmaz, J. Magalhães, P. Castells, N. Ferro, M. J. Silva, and F. Martins. Vol. 12036. Lecture Notes in Computer Science. Springer, 2020, pp. 449-453. doi: 10.1007/978-3-030-45442-5_56
- J** J. Devezas and S. Nunes. "Hypergraph-of-entity: A unified representation model for the retrieval of text and knowledge". In: Open Computer Science 9:1 (June 2019), pp. 103-127. doi: 10.1515/comp-2019-0006
- C** J. Devezas and S. Nunes. "Characterizing the Hypergraph-of-Entity Representation Model". In: Complex Networks and Their Applications VIII - Volume 2 Proceedings of the Eighth International Conference on Complex Networks and Their Applications COMPLEX NETWORKS 2019, Lisbon, Portugal, December 10-12, 2019, 2019, pp. 3-14. doi: 10.1007/978-3-030-36683-4_1
- C** J. Devezas and S. Nunes. "Graph-of-Entity: A Model for Combined Data Representation and Retrieval". In: Proceedings of the 8th Symposium on Languages, Applications and Technologies (SLATE 2019), Vila do Conde, Portugal, 2019. doi: 10.4230/OASlcs.SLATE.2019.1
- D** J. Devezas and S. Nunes. Simple English Wikipedia Link Graph with Clickstream Transitions 2018-12 [dataset]. INESC TEC research data repository, Mar. 2019. doi: 10.25747/83vk-zt74
- C** J. L. Devezas, S. Nunes, A. Guillén, Y. Gutiérrez, and R. Muñoz. "FEUP at TREC 2018 Common Core Track - Reranking for Diversity using Hypergraph-of-Entity and Document Profiling". In: Proceedings of the Twenty-Seventh Text REtrieval Conference, TREC 2018, Gaithersburg, Maryland, USA, November 14-16, 2018. Ed. by E. M. Voorhees and A. Ellis. Vol. 500-331. NIST Special Publication. National Institute of Standards and Technology (NIST), 2018. url: <https://trec.nist.gov/pubs/trec27/papers/FEUP-CC.pdf>
- C** J. Devezas and S. Nunes. "Social Media and Information Consumption Diversity". In: Proceedings of the Second International Workshop on Recent Trends in News Information Retrieval co-located with 40th European Conference on Information Retrieval (ECIR 2018), Grenoble, France, March 26, 2018. 2018, pp. 18-23. url: <http://ceur-ws.org/Vol-2079/paper5.pdf>
- C** J. L. Devezas, C. T. Lopes, and S. Nunes. "FEUP at TREC 2017 OpenSearch Track Graph-Based Models for Entity-Oriented". In: Proceedings of The Twenty-Sixth Text REtrieval Conference, TREC 2017, Gaithersburg, Maryland, USA, November 15-17, 2017. Ed. by E. M. Voorhees and A. Ellis. Vol. 500-324. NIST Special Publication. National Institute of Standards and Technology (NIST), 2017. url: <https://trec.nist.gov/pubs/trec26/papers/FEUP-O.pdf>
- R** J. Devezas. Army ANT: A Workbench for Innovation in Entity-Oriented Search - External Option: Scientific Activities - TREC Open Search. Research rep. Faculty of Engineering, University of Porto, June 2017. url: <https://hdl.handle.net/10216/110181>
- R** J. Devezas. Auditing Open Access Repositories - Free Option: Supervised Study - Digital Archives and Libraries. Research rep. Faculty of Engineering, University of Porto, May 2017. url: <https://hdl.handle.net/10216/104152>
- C** J. Devezas and S. Nunes. "Information Extraction for Event Ranking". In: 6th Symposium on Languages, Applications and Technologies (SLATE 2017). Ed. by R. Queirós, M. Pinto, A. Simões, J. P. Leal, and M. J. Varanda. Vol. 56. OpenAccess Series in Informatics (OASlcs). Dagstuhl, Germany: Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2017, 18:1-18:14. isbn: 978-3-95977-056-9. doi: 10.4230/OASlcs.SLATE.2017.18
- J** J. Devezas and S. Nunes. "Graph-Based Entity-Oriented Search: Imitating the Human Process of Seeking and Cross Referencing Information". In: ERCIM News. Special Issue: Digital Humanities 111 (Oct. 2017), pp. 13-14. url: <https://ercim-news.ercim.eu/en111/special/graph-based-entity-oriented-search-imitating-the-human-process-of-seeking-and-cross-referencing-information>
- C** T. Devezas, J. Devezas, and S. Nunes. "Exploring a Large News Collection Using Visualization Tools". In: Proceedings of the First International Workshop on Recent Trends in News Information Retrieval colocated with 38th European Conference on Information Retrieval (ECIR 2016), Padua, Italy, March 20, 2016. 2016, pp. 48-53. url: <http://ceur-ws.org/Vol-1568/paper9.pdf>
- C** J. L. Devezas and S. Nunes. "Index-Based Semantic Tagging for Efficient Query Interpretation". In: Experimental IR Meets Multilinguality, Multimodality, and Interaction - 7th International Conference of the CLEF Association, CLEF 2016, Évora, Portugal, September 5-8, 2016, Proceedings. Ed. by N. Fuhr, P. Quaresma, T. Gonçalves, B. Larsen, K. Balog, C. Macdonald, L. Cappellato, and N. Ferro. Vol. 9822. Lecture Notes in Computer Science. Springer, 2016, pp. 208-213. doi: 10.1007/978-3-319-44564-9_17

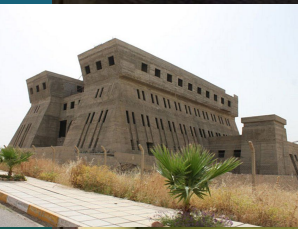
Appendix

Extra detail to aid discussion

Historical perspective

For centuries, information has been organized, stored, and retrieved

- Clay tablets in Ashurbanipal
- Books in modern libraries
- Digitally encoded documents in computers
- Entities and their relations in knowledge bases



Library of Ashurbanipal (The British Museum)



Chief Librarian

As the library grew, Ashurbanipal appointed a chief librarian to oversee the collection. This role was crucial in maintaining the integrity and accessibility of the vast archive of knowledge.

Acquisition

The library's collection was not static; it grew through the acquisition of new tablets and fragments. These were often obtained through trade or as spoils of war, reflecting the library's role as a center of knowledge and power.

Restoration

Many of the tablets in the library were damaged or incomplete. The library's staff worked to restore these fragments, piecing together the original text and preserving the knowledge for future generations.

Preservation

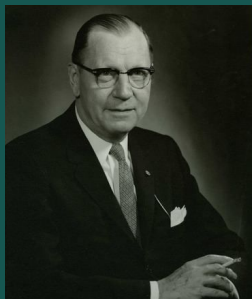
The library's collection was preserved through a variety of methods, including the use of clay tablets and the creation of duplicate copies. This ensured that the knowledge stored in the library was not lost to time.

Dissemination

The library's collection was not just a storehouse of knowledge; it was a place where knowledge was shared and disseminated. The library's staff made the tablets available to scholars and students, ensuring that the knowledge was passed on.

Reconstruction

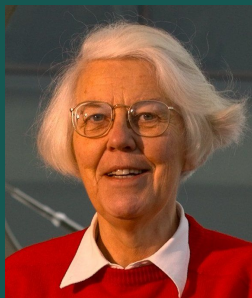
The library's collection was reconstructed through the use of clay tablets and the creation of duplicate copies. This ensured that the knowledge stored in the library was not lost to time.



Term Frequency

The weight of a term that occurs in a document is simply proportional to the term frequency.

– Hans Peter Luhn, 1957



Inverse Document
Frequency

The specificity of a term can be quantified as an inverse function of the number of documents in which it occurs.

– Karen Spärck Jones, 1972



Inverted Files

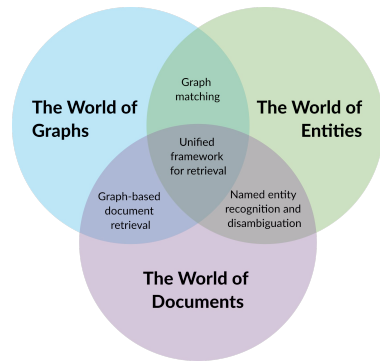
The first important class of techniques for secondary key retrieval is based on the idea of an inverted file. This does not mean that the file is turned upside down; it means that the roles of records and attributes are reversed. Instead of listing the attributes of a given record, we list the records having a given attribute.

– Donald Ervin Knuth, 1973

Consolidating models

- From physics to machine learning
 - Efforts to unify theories and models
- Towards general approaches to IR
 - Identifying commonalities along the pipeline and tasks
- Graphs as general representation models
 - Combining all available information sources
- Unified framework for IR
 - Solving the information need is the only task

Three axes covered



Classical models

- Virtual documents
- Triplestores
- Combined signals (single task or chained tasks)
- Joint indexing of text and triples (very few contributions)

Learning to rank

- Entity profiles represented as virtual documents
- Entity features
- Joint learning of word and entity representations

Graph-based models

- Link analysis
- Text as a graph
- Knowledge graphs
- Entity graph from text
- Entity graph as a tensor
- Graph matching
- Hypergraph-based
- Random walk based

Classical models

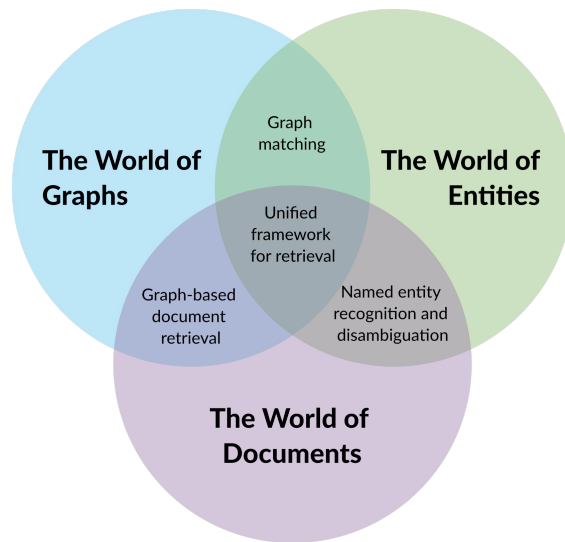
- Ranking
 - From TF-IDF
 - To Markov networks
- Representation
 - Virtual documents
 - Triplestores
- Hybrid approaches
 - Combined signals (single task or chained tasks)
 - Joint indexing of text and triples (very few contributions)

Learning-to-rank models

- Entity ranking based on:
 - Entity features
 - Text features from the Wikipedia article
 - Graph features from the knowledge graph
 - Entity profiles represented as virtual documents
 - “Flattened” data from RDF
 - Passages of text mentioning the entity
- Joint learning of representations for words and entities

Graph-based models

- Link analysis
- Text as a graph
- Knowledge graphs
- Entity graph from text
- Entity graph as a tensor
- Graph matching
- Hypergraph-based
- Random walk based



Anchor / core references

Title:

Entity-Oriented Search

Authors:

K. Balog

Year:

2018

DOI:

10.1007/978-3-319-93935-3

- First complete reference in the area
- Clear definitions of fundamental concepts
- Identification of tasks and applications
- Provides a compilation and a convergence of information

Anchor / core references

Title:

Concordance-Based Entity-Oriented Search

Authors:

M. Bautin and S. Skiena

Year:

2007

DOI:

10.1109/WI.2007.84

- Analyzes the presence of entities in queries (based on AOL query log):
 - 18-39% queries directly reference entities
 - 73-87% queries contain at least one entity
- “First-in-literature” implementation of an entity search engine
- Archetype for approaches based on virtual documents

Anchor / core references

Title:

Ad-hoc object retrieval in the web of data

Authors:

J. Pound, P. Mika, and H. Zaragoza

Year:

2010

DOI:

10.1145/1772690.1772769

Five query categories for ad hoc entity retrieval:

- Entity query
- Type query
- Attribute query
- Relation query
- Keyword query



Applied to the ANT search engine

Anchor / core references

Title:

An Index for Efficient Semantic Full-text Search

Authors:

H. Bast and B. Buchhold

Year:

2013

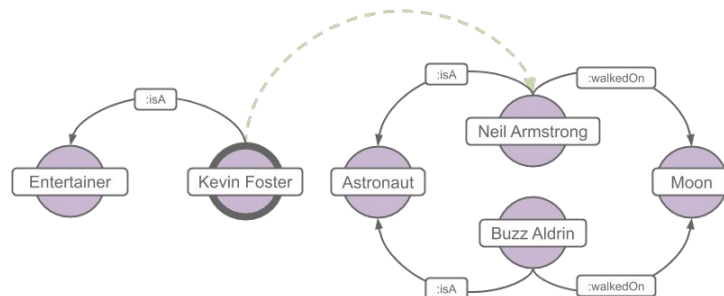
DOI:

10.1145/2505515.2505689

Query: entertainers that are friends with astronauts who walked on the moon

"After the act, Kevin Foster went down to the audience, to hug his friend, Neil Armstrong, who had been sitting in the crowd since the beginning of the show."

Neil Armstrong	:isA	Astronaut
Neil Armstrong	:walkedOn	Moon
Buzz Aldrin	:isA	Astronaut
Buzz Aldrin	:walkedOn	Moon
Kevin Foster	:isA	Entertainer



Anchor / core references

Title:

Graph-of-word and TW-IDF: new approach to ad hoc IR

Authors:

F. Rousseau and M. Vazirgiannis

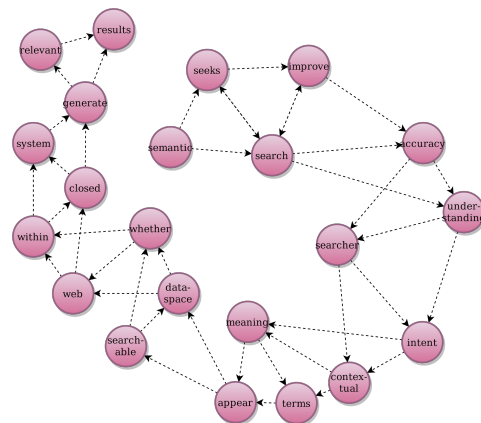
Year:

2013

DOI:

10.1145/2505515.2505671

- Nodes represent terms
- Edges represent following terms within a window of size N
- TW is given by the indegree



Anchor / core references

Title:

Modeling Higher-order Term
Dependencies in Information Retrieval
Using Query Hypergraphs

Authors:

M. Bendersky and W. B. Croft

Year:

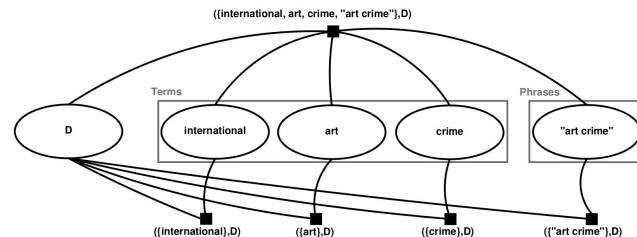
2012

DOI:

10.1145/2348283.2348408

Query hypergraph model

- Log-linear retrieval model
- Solved through a factor graph
- Similar to Markov networks
- But captures higher-order relations (e.g., bigrams, named entities)



Anchor / core references

Title:

It's more than just overlap: Text As Graph

Authors:

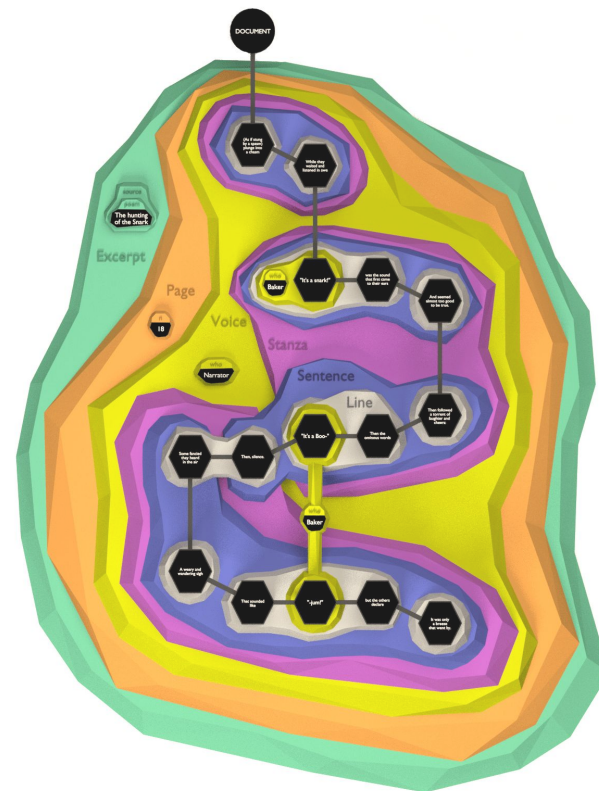
R. Haentjens Dekker and D. J. Birnbaum

Year:

2017

DOI:

10.4242/BalisageVol19.Dekker01



COMBINED DATA

Combined data is a collection of corpora and knowledge bases, which includes not only the natural relations between documents (e.g., hyperlinks in the web), and entities (e.g., object properties in triplestores), but also cross-context relations, from mentions found in documents to entities in knowledge bases, and from entities found in knowledge bases to instances of the same entity in other knowledge bases.

THESIS STATEMENT

A graph-based joint representation of unstructured and structured data has the potential to unlock novel ranking strategies, that are, in turn, able to support the generalization of entity-oriented search tasks and to improve overall retrieval effectiveness by incorporating explicit and implicit information derived from the relations between text found in corpora and entities found in knowledge bases.

THESIS STATEMENT



A graph-based joint representation of unstructured and structured data has the potential to unlock novel ranking strategies, that are, in turn, able to support the generalization of entity-oriented search tasks and to improve overall retrieval effectiveness by incorporating explicit and implicit information derived from the relations between text found in corpora and entities found in knowledge bases.

Test collections

INEX 2009 Wikipedia collection

- 2.6M XML documents
- 2010 Ad Hoc track
 - Document ranking
- 2009 XER track
 - Entity ranking
 - List completion

TREC Washington Post Corpus

- 595K JSON documents
- 2018 Common Core track

Social Science Open Access Repository

- 32K JSON documents
- 2017 OpenSearch track
- Team-draft interleaving

Contributed datasets

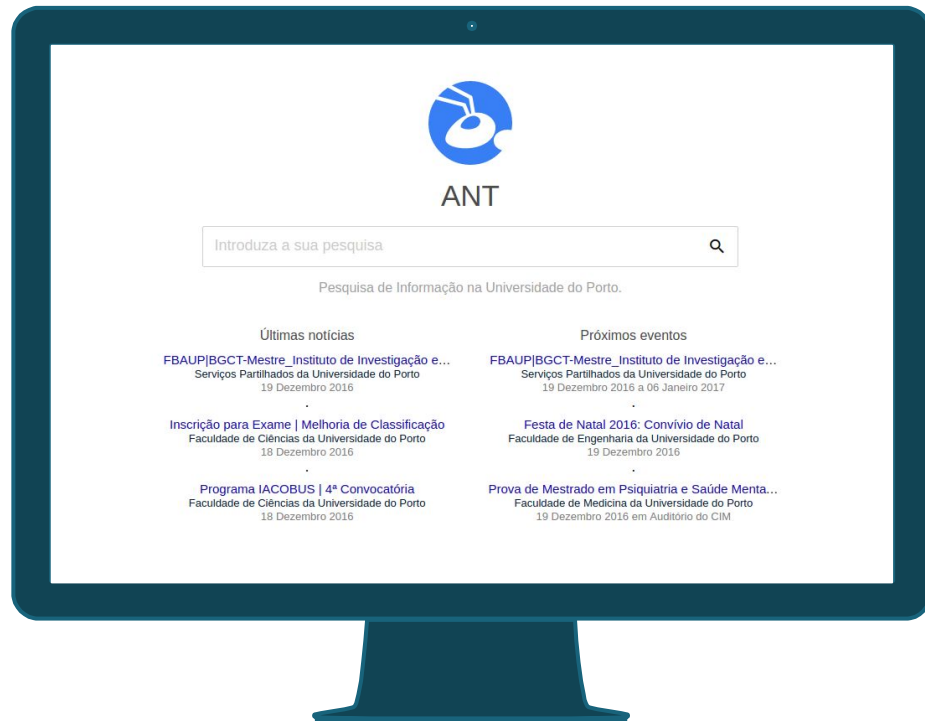
Simple English Wikipedia Link Graph with Clickstream
Transitions 2018-12

DOI: 10.25747/83vk-zt74

ANT

Entity-oriented search engine for the University of Porto

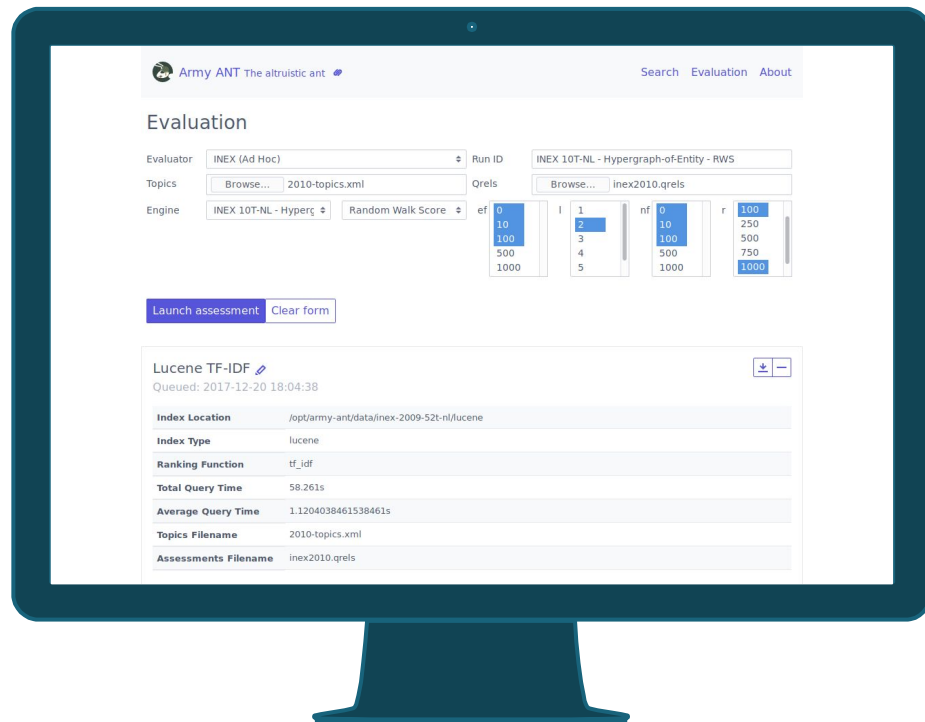
- Working prototype (<https://ant.fe.up.pt>)
- Exposure to ~1,000 weekly users
- Manifested interest by some of the faculty's content managers



Army ANT

Workbench for innovation in
entity-oriented search

- Promotes freedom and exploration
- Supports IR research in a flexible way
- Available at GitHub
(<https://github.com/feup-infolab/army-ant>)



Hypergraph-of-entity: characterization

Basic statistics over time (i.e., as the index grows):

- Number of nodes and hyperedges per type and direction
- Node-based and hyperedge-based degree distributions
- Hyperedge cardinality distribution
- Clustering coefficient
- Average path length & diameter
- Density

- Shortest distances computed based on random walks
- Two-node clustering coefficients
 - Based on a set of sampled nodes
 - And a large sample of their neighbors
- Density based on a corresponding bipartite graph
 - Hyperedge-cardinality notation recognized as useful by the community

$$D = \frac{2 \sum_k k |E_U^k| + \sum_{k_1, k_2} (k_1 + k_2) |E_D^{k_1, k_2}|}{2(n + m)(n + m - 1)}$$

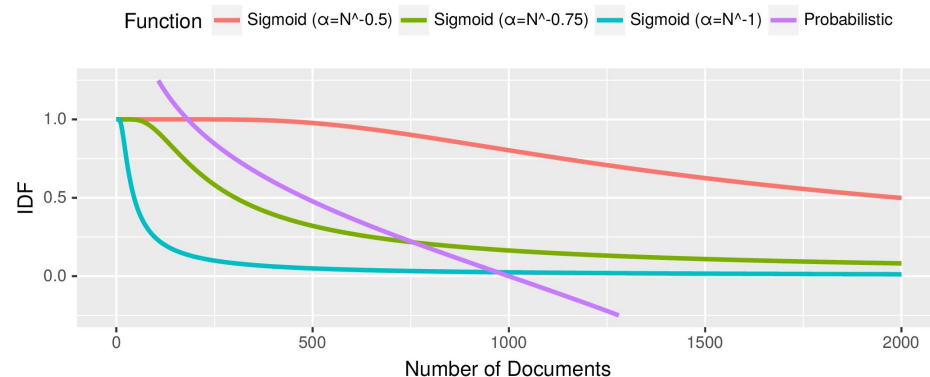
Hypergraph-of-entity: joint representation model

Extensions: weights

Table 7.2: Hypergraph-of-entity weighting functions.

(a) Nodes.

Node / Weight	Description
term $2S\left(\alpha \frac{N - n_t}{n_t}\right) - 1$	<p>We used a variation of the IDF, with a tunable $\alpha < 1$ parameter to control how fast the function decreases.</p> <ul style="list-style-type: none"> - S is the sigmoid function - N is the number of documents in the collection - n_t is the number of documents where a given term t occurs. - We used $\alpha = N^{-0.75}$.
entity <p>Same as <i>term</i>.</p>	<p>In the future, we will experiment with different values of α for terms and entities, in particular alternative exponents to -0.75.</p>



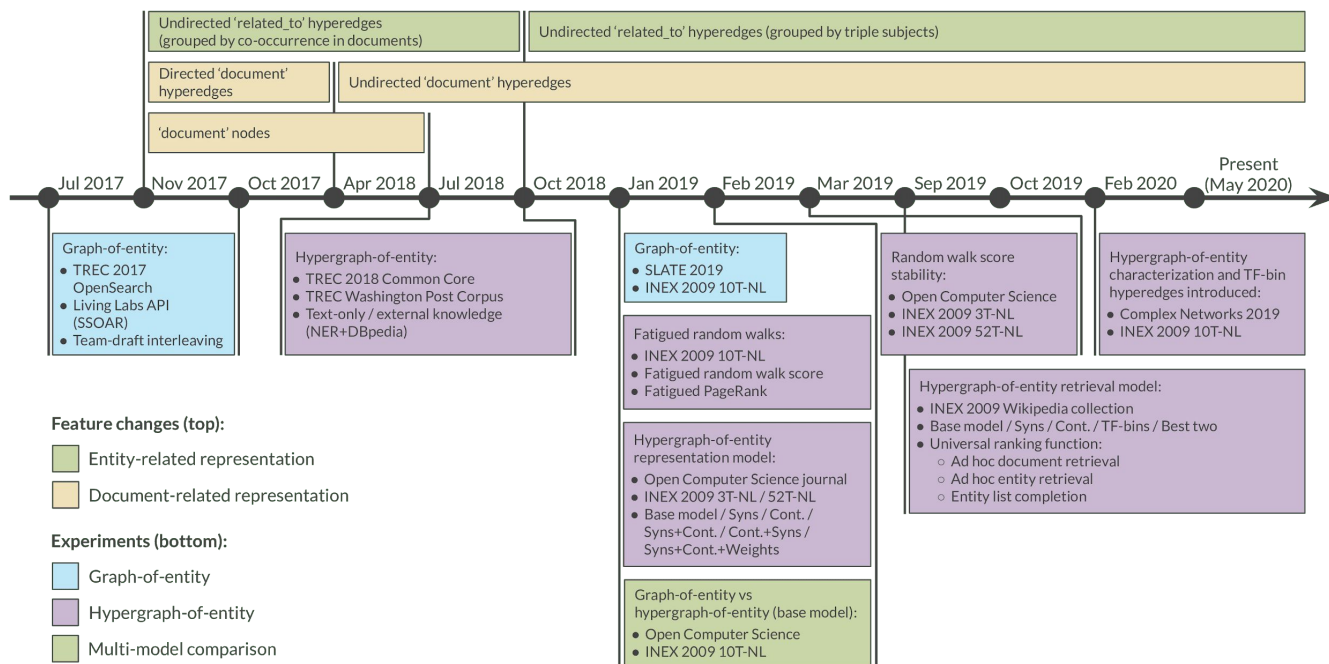
Hypergraph-of-entity: joint representation model

Extensions: weights

(b) Hyperedges.

Hyperedge / Weight	Description
<i>document</i> 0.5	Linking a term or entity simply through document co-occurrence is weak, so we use a constant weight lower than one.
<i>related_to</i> $\frac{1}{ e_r } \sum_{v \in e_r} \frac{ \{u \in e'_r : e'_r \in E_r \setminus \{e_r\} \wedge v \in e'_r\} }{ e_r }$	For each entity within the hyperedge, we calculate the fraction of reachable other entities and average all results. - E_r is the set of all <i>related_to</i> hyperedges. - $e_r \in E_r$ is the specific <i>related_to</i> hyperedge, for which we are calculating the weight.
<i>contained_in</i> $\frac{1}{ t }$	Links with fewer terms t , where t refers to the tail set in $(t, h) \in E_c \wedge t \subseteq V_t$, should be more frequently followed, since the certainty that the hyperedge leads to the entity is higher.
<i>synonym</i> $\frac{1}{ e_s }$	The higher the number of possible synonyms $e_s \in E_s \wedge e_s \subseteq V_t$, the less certain we are about the hyperedge — we rely on the synonyms of the first (and most probable) sense according to WordNet.
<i>context</i> $\frac{1}{ e_x } \sum_{t_i \in e_x \setminus \{t_k\}} \frac{\text{sim}(t_k, t_i) - \min_{\text{sim}}}{1 - \min_{\text{sim}}}$	A context $e_x \in E_x$ is only as good as the average of all similarities between the original term $t_k \in e_x$ and all other terms $t_i \in e_x \setminus \{t_k\}$. We normalize the weight taking into account the threshold used to create the word2vec SimNet.

Evaluation: experimentation timeline



Evaluation: main experiments (pt. 2)

Rank stability:

- Average Kendall's coefficient of concordance
 - Over 100 similar runs per configuration
 - For different values of ℓ and r
- For $\ell \in \{2, 3, 4\}$:
 - 84-90% stable for $r=100$
 - 94-97% stable for $r=1,000$
 - 99% stable for $r=10,000$

Evaluation: INEX 2009 Wikipedia collection

Three subsets:

- INEX 2009 3T-NL (2.2k docs)
- INEX 2009 10T-NL (7.5k docs)
- INEX 2009 52T-NL (37.8k docs)

Created through:

- Random sampling of n topics
- Retained relevance judgments for selected topics
- Retained only judged documents

However:

- Exclusively for assessing ad hoc document retrieval
- Based on the qrels for the 2010 Ad Hoc track
- Goal: eventually index the full collection
- Challenge: scalability

Discussion

Ad hoc document retrieval

- Best MAP: **0.1689** (vs 0.1710 TF-IDF)
 - Base model
 - INEX 2009 10T-NL
 - $\ell=2$, $r=10,000$, $\text{exp.}=\text{true}$
- Best P@10: **0.2692** (vs 0.0692 TF-IDF)
 - Synonyms+Context model
 - INEX 2009 full collection
 - $\ell=2$, $r=10,000$, $\text{exp.}=\text{false}$

Note: In bold the best scores for the hypergraph-of-entity; in parenthesis the baseline result of Lucene TF-IDF for the same experiment.

Discussion

Ad hoc entity retrieval

- Best MAP: **0.1390** (vs 0.0373 TF-IDF)
 - Base model
 - INEX 2009 full collection
 - $\ell=2$, $r=10,000$, $\text{exp.}=\text{false}$
- Best P@10: **0.2509** (vs 0.0636 TF-IDF)
 - Synonyms+Context model
 - INEX 2009 full collection
 - $\ell=2$, $r=10,000$, $\text{exp.}=\text{false}$

Entity list completion

- Best MAP: **0.0884** (vs 0.0558 TF-IDF)
 - Synonyms+Context model
 - INEX 2009 full collection
 - $\ell=2$, $r=10,000$, $\text{exp.}=\text{false}$
- Best P@10: **0.0788** (vs 0.1000 TF-IDF)
 - Synonyms+Context model
 - INEX 2009 full collection
 - $\ell=2$, $r=10,000$, $\text{exp.}=\text{false}$

Note: In bold the best scores for the hypergraph-of-entity; in parenthesis the baseline result of Lucene TF-IDF for the same experiment.