

# **ENTITY-BASED QUERY INTERPRETATION**

BACHELOR'S DEFENCE

MARCEL GOHSEN

BAUHAUS-UNIVERSITÄT WEIMAR

04 JULY 2018

# PROBLEM OF QUERY INTERPRETATION

new york times square dance

# PROBLEM OF QUERY INTERPRETATION

new york times

square dance

"All the News  
That's Fit to Print"

# The New York Times

VOL. CLVIII , No. 54,678      © 2009 The New York Times      NEW YORK, SUNDAY, MAY 17, 2009      \$5 beyond the greater New York metropolitan area. \$4.00



## Square Dance

19 U.S. States Have Designated It As Their Official State Dance

**Late Edition**  
Today's shower, picking to 80°, high 60°, low 50°, winds 10 mph, to 60°, sunny and cool, high 64°. Yesterday's high, 70°, low, 54°. Weather map and details, Page 24.

### From a Theory To a Consensus On Emissions

Permits Gain Political Edge Over Taxation

By JOHN M. BROOKER

WASHINGTON — As Congress weighs imposing a mandatory limit on greenhouse-gas emissions — an outcome still far from certain — it is likely to turn to a system of permits, a market-based ceiling on total emissions and allowances polling industries to buy and sell permits.

That approach, known as cap and trade, has been embraced by President Obama, Democratic leaders in Congress and environmental groups and a growing number of business interests. It offers a way to combat rising industries like autos, steel and aluminum.

But just a decade ago, many of today's supporters dismissed the idea of tradable emissions permits as an inefficient way for politicians to meet the real costs of cutting air pollution. The right answer, they said, was strict government regulation, state-of-the-art technology and a federal tax on every ton of harmful emissions.

Now, though, many of those who once opposed the idea have come around. Last year, for example, most of the Senate's supporters dismissed the idea of tradable emissions permits as an inefficient way for politicians to meet the real costs of cutting air pollution. The right answer, they said, was strict government regulation, state-of-the-art technology and a federal tax on every ton of harmful emissions.

And if he nominates Kathleen M. Sullivan, a law professor at Stanford, they plan to denounce her as a "rogue" supporter of gay marriage.

If he nominates Judge Sonia Sotomayor, they plan to accuse her of being "willing to expand constitutional rights beyond the boundaries of common sense."

Prepared to impose the confirmation of Justice Sonia Sotomayor, the Senate Republicans are gathering to stockpile ammunition. Ten memorandums summarizing the conservative positions of the New York Times, provide a window onto how they hope to frame the confirmation debate.

The memorandums dissect possible nominees' records, noting statements the groups find objectionable on issues like abortion.

### CONSERVATIVES MAP STRATEGIES ON COURT FIGHT

### MEMOS OUTLINE ATTACKS

Hoping to Re-Energize G.O.P. by Opposing Obama's Choice

By CHARLIE SAVAGE

WASHINGTON — If President Obama nominates Solicitor General Elena K. Kagan to replace Justice Ruth Bader Ginsburg on the Supreme Court, conservatives plan to attack her as an "outspoken" supporter of gay marriage and abortion rights.

If he nominates Judge Sonia Sotomayor, they plan to accuse her of being "willing to expand constitutional rights beyond the boundaries of common sense."

And if he nominates Kathleen M. Sullivan, a law professor at Stanford, they plan to denounce her as a "rogue" supporter of gay marriage.

Prepared to impose the confirmation of Justice Sonia Sotomayor, the Senate Republicans are gathering to stockpile ammunition. Ten memorandums summarizing the conservative positions of the New York Times, provide a window onto how they hope to frame the confirmation debate.

The memorandums dissect possible nominees' records, noting statements the groups find objectionable on issues like abortion.

# PROBLEM OF QUERY INTERPRETATION

new york

times square

dance



# ENTITIES IN QUERIES

## ■ Named Entity

- ▶ object from the real world with a proper name
- ▶ e.g., *person, location, organization*

## ■ Entities in Queries

- ▶ Definitions differ
- ▶ May be limited to proper nouns <sup>1</sup>
- ▶ May include general concepts <sup>2</sup>

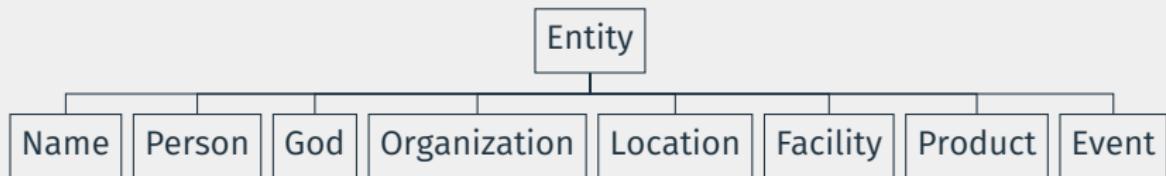
---

<sup>1</sup>[Hasibi et al., 2015]

<sup>2</sup>[Cornolti et al., 2016]

# USED ENTITY TAXONOMY

- Based on “Extended Named Entity Hierarchy”  
[Sekine et al., 2002]
- 8 main classes
- 108 specialized subclasses



- for example: removed class *Units* (e.g., kilogram)

# **TRADITIONAL PROBLEM STATEMENTS**

## ENTITY LINKING [HASIBI ET AL., 2015]

Linking an entity in a query to the most likely candidate in some knowledge base.

obama mother → (“obama”, Barack Obama)

new york pizza manhattan → (“new york”, New York City)  
 (“manhattan”, Manhattan)

Issues:

- Non-overlapping entities only

## INTERPRETATION FINDING [HASIBI ET AL., 2015]

Finding subsets of semantic compatible non-overlapping linked entities

obama mother → {Barack Obama}

new york pizza manhattan → {New York City, Manhattan}

{New York-Style Pizza, Manhattan}

Issues:

- Imprecise interpretations
- Explicit mentioned entities only

## INTERPRETATION FINDING [HASIBI ET AL., 2015]

Finding subsets of semantic compatible non-overlapping linked entities

obama **mother** → {Barack Obama} **mother?**

new york **pizza** manhattan → {New York City, Manhattan} **pizza?**  
{New York-Style Pizza, Manhattan}

Issues:

- Imprecise interpretations
- Explicit mentioned entities only

# **REDEFINED PROBLEMS**

## EXPLICIT ENTITY RECOGNITION

**Given:** - Query

**Task:** - Identifying explicit mentioned entities in a query  
- Segment is an entity's name or surface form

obama mother → (“obama”, Barack Obama)

(“obama”, Michelle Obama)

(“obama”, Natsuki Obama)...

new york pizza manhattan → ("new york", New York City)

(“new york”, New York (state))

(“manhattan”, Manhattan

(“manhattan”, Manhattan (film))...

# IMPLICIT ENTITY RECOGNITION

## Given: - Query

**Task:** - Identifying implicitly referenced entities in a query  
- Segment is a description of an entity

new york pizza manhattan → Ø

president of usa → (“president of usa”, Donald Trump)  
("president of usa", Barack Obama)  
("president of usa", George W. Bush)

# ENTITY-BASED QUERY INTERPRETATION

**Given:**

- Query
- Explicit entities in query
- Implicit entities in query

**Task:**

- Semantically segmentation of query
- Replacing explicit and implicit entity-mentions with entities

obama mother → {Barack Obama, Ann Dunham}  
{Michelle Obama, Marian Shields}

...

new york pizza manhattan → {New York City, “pizza”, Manhattan}

...

# CORPORA

- Dataset of the ERD'14 Challenge
- 91 queries
  - ▶ 45 queries having annotated entities
- Provides query interpretation

obama family tree → {Barack Obama}

east ridge high school → {East Ridge High School (FL)}  
{East Ridge High School (MN)}  
{East Ridge High School (KY)}

## YSQL DATASET [YAHOO, 2010]

- “Yahoo Search Query Log to Entities”
- 2635 queries
  - ▶ 2583 queries having annotated entities
- No query interpretations

france 1998 final → France National Football Team,

France, Fifa World Cup 1998 Final

obama mother → Barack Obama, Ann Dunham

# DBPEDIA-ENTITY V2 DATASET [HASIBI ET AL., 2017]

- Collection for Entity Search
- 467 queries
- No query interpretations
- Introduced relevance levels
  - ▶ 2: highly relevant
  - ▶ 1: relevant
  - ▶ 0: irrelevant

john lennon, parents → {Julia Lennon : 2,  
Alfred Lennon : 1  
... : 0}

# QUERY INTERPRETATION CORPUS

- Queries from the three existing corpora
- Manually (re-)annotated:
  - ▶ Query difficulty judgments {easy | moderate | hard}
  - ▶ Explicit entities with relevance judgments {relevant | plausible}
  - ▶ Implicit entities with relevance judgments
  - ▶ Entity-based query interpretations with relevance judgments
- 2068 queries
  - ▶ 1578 queries with explicit entities
  - ▶ 131 queries with implicit entities
  - ▶ 1597 queries with query interpretations

# **ALGORITHMIC APPROACHES**

# ENTITY LINKING STEPS

Typical steps for entity linking frameworks

- (i) **Candidate Generation**
- (ii) **Scoring**
- (iii) **Selecting**

## (I) CANDIDATE GENERATION

- DBpedia Ontology [DBpedia, 2017] used for classification
  - ▶ Digital representation of our entity taxonomy
- Index all Wikipedia articles that represent entities
- Retrieve the top 100 articles from the index containing a segment from the query
- Retrieve for each segment of the query

## (II) SCORING

$$Jaccard(T_1, T_2) = \frac{|T_1 \cap T_2|}{|T_1 \cup T_2|}$$

$$norm = \frac{|\textit{segment}|}{|\textit{query}|}$$

### (III) SELECTION

- Precision vs. Recall
- Threshold vs. Fixed number of retrieved entities
- Take the top 20 entities by score

# EVALUATION

# EVALUATION RESULTS FOR EXPLICIT ENTITY RECOGNITION

<b>Algorithm</b>	<i>rec</i>	<i>prec</i>	$F_1$	<i>rec*</i>	$F_1^*$	<i>RT</i>
Nordlys EL	.55	.69	.58	.50	.52	4400 ms
Explicit Entity Approach	.40	.16	.18	.35	.16	270 ms
Smaph	.38	.45	.37	.32	.31	117000 ms
TagMe	.37	.39	.33	.31	.28	<b>40 ms</b>
Nordlys ER	.33	.05	.07	.29	.06	1900 ms
Baseline	.26	.26	.26	.26	.26	-

# CONCLUSION

- Refined problem statements for entity linking
  - ▶ Ambiguous explicit and implicit entities
  - ▶ More precise and diverse query interpretations
- Query Interpretation Corpus
  - ▶ Comparatively large corpus
  - ▶ Explicit and implicit entities
  - ▶ Query interpretations
- Algorithmic Approaches
  - ▶ Efficient explicit entity recognition
  - ▶ Implicit entity recognition prototype

Thank you for the attention!

# REFERENCES I

 CARMEL, D., CHANG, M.-W., GABRILOVICH, E., HSU, B.-J. P., AND WANG, K. (2014).

**ERD'14: ENTITY RECOGNITION AND DISAMBIGUATION CHALLENGE.**  
*SIGIR Forum*, 48(2):63–77.

 CORNOLTI, M., FERRAGINA, P., CIARAMITA, M., RÜD, S., AND SCHÜTZE, H. (2016).

**A PIGGYBACK SYSTEM FOR JOINT ENTITY MENTION DETECTION AND LINKING IN WEB QUERIES.**

In *Proceedings of the 25th International Conference on World Wide Web*, WWW '16, pages 567–578, Republic and Canton of Geneva, Switzerland. International World Wide Web Conferences Steering Committee.

 DBPEDIA (2017).

**DBPEDIA ONTOLOGY 2016-10.**

<https://wiki.dbpedia.org/services-resources/ontology>.

## REFERENCES II

-  HASIBI, F., BALOG, K., AND BRATSBERG, S. E. (2015).  
**ENTITY LINKING IN QUERIES: TASKS AND EVALUATION.**  
In Allan, J., Croft, W. B., de Vries, A. P., and Zhai, C., editors, *Proceedings of the 2015 International Conference on The Theory of Information Retrieval, ICTIR 2015, Northampton, Massachusetts, USA, September 27-30, 2015*, pages 171–180. ACM.
-  HASIBI, F., NIKOLAEV, F., XIONG, C., BALOG, K., BRATSBERG, S. E., KOTOV, A., AND CALLAN, J. (2017).  
**DBPEDIA-ENTITY V2: A TEST COLLECTION FOR ENTITY SEARCH.**  
In Kando, N., Sakai, T., Joho, H., Li, H., de Vries, A. P., and White, R. W., editors, *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval, Shinjuku, Tokyo, Japan, August 7-11, 2017*, pages 1265–1268. ACM.
-  SEKINE, S., SUDO, K., AND NOBATA, C. (2002).  
**EXTENDED NAMED ENTITY HIERARCHY.**  
In LREC.

## REFERENCES III



YAHOO (2010).

**L24 - YAHOO SEARCH QUERY LOG TO ENTITIES V1.0.**

<https://webscope.sandbox.yahoo.com/>.





# EVALUATION METRICS

$$prec = \begin{cases} \frac{|E \cap E'|}{|E|}, & \text{if } |E| > 0 \\ 1, & \text{if } |E| = 0, |E'| = 0 \\ 0, & \text{if } |E| = 0, |E'| > 0 \end{cases} \quad (1)$$

$$rec = \begin{cases} \frac{|E \cap E'|}{|E'|}, & \text{if } |E'| > 0 \\ 1, & \text{if } |E| = 0, |E'| = 0 \\ 0, & \text{if } |E| > 0, |E'| = 0 \end{cases} \quad (2)$$

$$F_1 = \frac{2 \cdot prec \cdot rec}{prec + rec} \quad (3)$$

# EVALUATION METRICS

$$w = \frac{\sum_{e \in E \cap E'} \text{rel}(e)}{\sum_{e' \in E'} \text{rel}(e')} \quad (4)$$

$$\text{rec}^* = w \cdot \text{rec} \quad (5)$$

$$F_1^* = \frac{2 \cdot \text{prec} \cdot \text{rec}^*}{\text{prec} + \text{rec}^*} \quad (6)$$

<b>Algorithm</b>	<i>prec</i>	<i>rec</i>	$F_1$	<i>rec*</i>	$F_1^*$
TagMe	.52	.49	.44	.42	.37
Smaph	.58	.48	.47	.40	.39
Explicit Entity Approach	.14	.47	.17	.40	.14
Nordlys EL	.64	.45	.49	.38	.41
Nordlys ER	.04	.43	.07	.37	.07