

# SEMANTIC SEARCH WITH KEYWORD QUERIES

Master's Thesis

Eugen Sawin

*Chair of Algorithms and Data Structures  
University of Freiburg*

# FULL-TEXT SEARCH

User Query

Films directed by Stanley Kubrick

# FULL-TEXT SEARCH

## User Query

Films directed by Stanley Kubrick

## Results

- ▶ **Stanley Kubrick - IMDb**  
*[www.imdb.com/name/nm0000040/](http://www.imdb.com/name/nm0000040/)*
- ▶ **Stanley Kubrick - Wikipedia**  
*[en.wikipedia.org/wiki/Stanley\\_Kubrick](http://en.wikipedia.org/wiki/Stanley_Kubrick)*
- ▶ **Stanley Kubrick, Film Director Dies at 70**  
*[www.nytimes.com/.../movies/stanley-kubrick...](http://www.nytimes.com/.../movies/stanley-kubrick...)*

# FULL-TEXT SEARCH

## User Query

Films directed by Stanley Kubrick

## Results

- ▶ Stanley Kubrick - IMDb  
[www.imdb.com/name/nm0000040/](http://www.imdb.com/name/nm0000040/)
- ▶ Stanley Kubrick - Wikipedia  
[en.wikipedia.org/wiki/Stanley\\_Kubrick](http://en.wikipedia.org/wiki/Stanley_Kubrick)
- ▶ Stanley Kubrick, Film Director Dies at 70  
[www.nytimes.com/.../movies/stanley-kubrick...](http://www.nytimes.com/.../movies/stanley-kubrick...)

We asked for **films** – got **documents**

# SEMANTIC SEARCH

User Query

Films directed by Stanley Kubrick

# SEMANTIC SEARCH

## User Query

Films directed by Stanley Kubrick

## Results

- ▶ A Clockwork Orange
- ▶ 2001: A Space Odyssey
- ▶ Dr. Strangelove or How I Learned...

# SEMANTIC SEARCH

## User Query

Films directed by Stanley Kubrick

## Results

- ▶ A Clockwork Orange
- ▶ 2001: A Space Odyssey
- ▶ Dr. Strangelove or How I Learned...

We asked for films – got film entities

# WHY SEMANTIC SEARCH?

- ▶ Over 40% of web searches are **entity searches**
- ▶ Focused results **save time**
- ▶ Suitable for machine consumption and **voice output**
- ▶ *(Finds results where document retrieval fails)*



# WHY SEMANTIC SEARCH?

- ▶ Over 40% of web searches are **entity searches**
- ▶ Focused results **save time**
- ▶ Suitable for machine consumption and **voice output**
- ▶ *(Finds results where document retrieval fails)*

Evolution of **intelligent search**

# WHY KEYWORD QUERIES?

# WHY KEYWORD QUERIES?

- ▶ Simple interface
- ▶ No expert knowledge required
  - ▶ Query languages
  - ▶ System-imposed limitations
- ▶ Effective for both text and voice input
- ▶ Users don't need to adapt

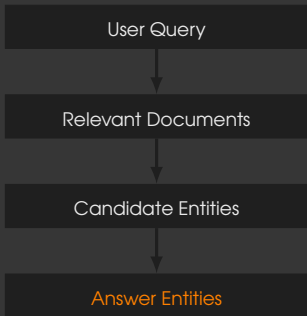
# WHY KEYWORD QUERIES?

- ▶ Simple interface
- ▶ No expert knowledge required
  - ▶ Query languages
  - ▶ System-imposed limitations
- ▶ Effective for both text and voice input
- ▶ Users don't need to adapt

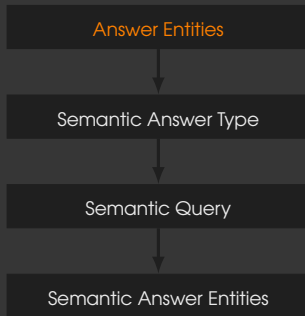
Semantic search for human beings

# TWO-PHASE APPROACH

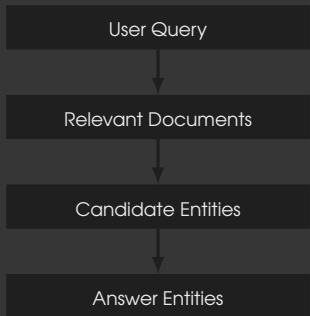
## Entity Retrieval



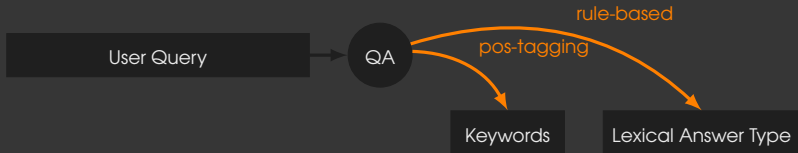
## Deep Search



# ENTITY RETRIEVAL PHASE



# QUERY ANALYSIS



- ▶ **Keywords:** nouns, adjectives and verbs
- ▶ **LAT:** first non-possessive noun

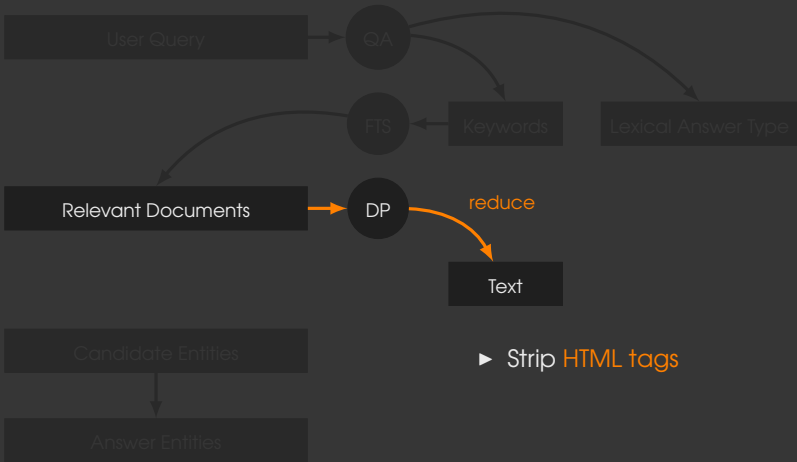
# DOCUMENT RETRIEVAL



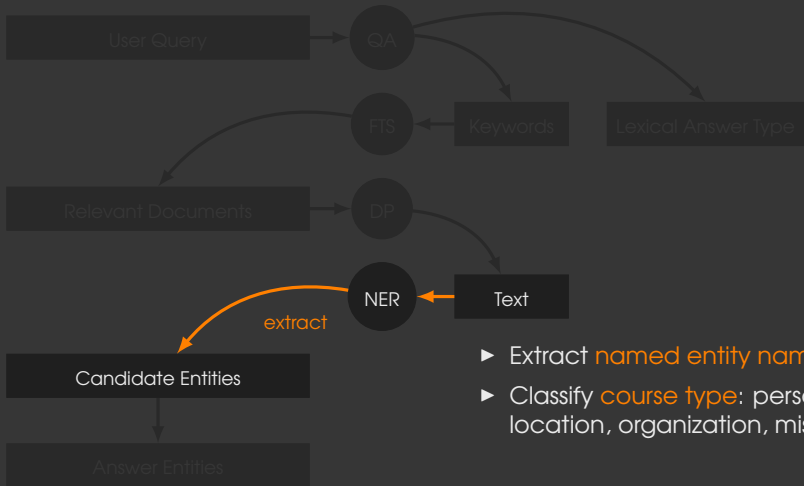
▶ Google Custom Search API:  
top 10 results only



# DOCUMENT SEGMENTATION

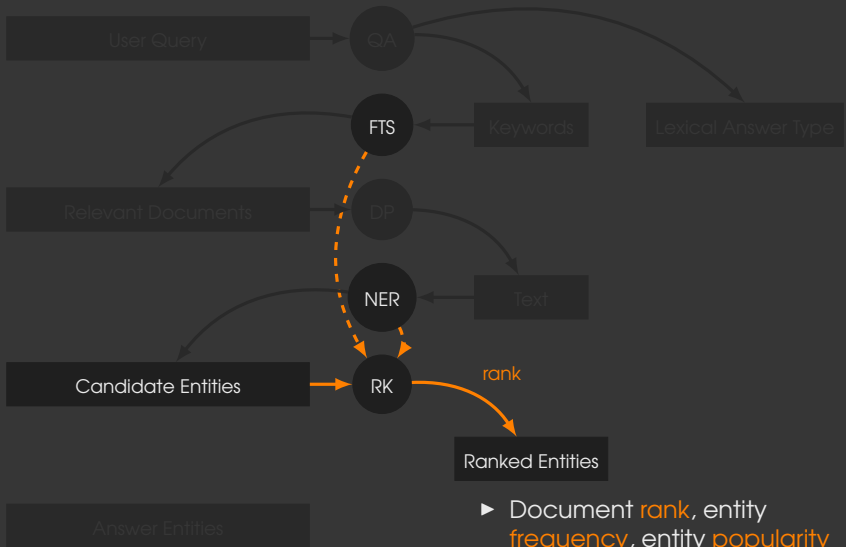


# ENTITY EXTRACTION



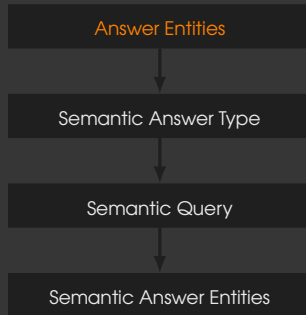
- ▶ Extract **named entity names**
- ▶ Classify **course type**: person, location, organization, misc

# ENTITY RANKING

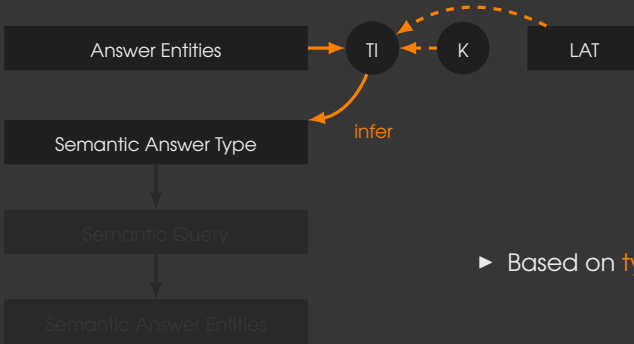




# DEEP SEARCH PHASE

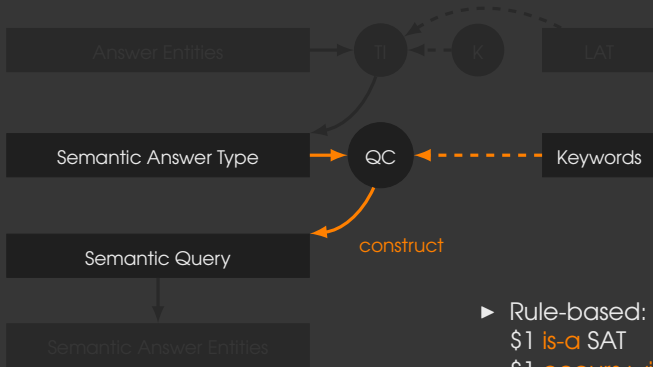


# TYPE INFERENCE



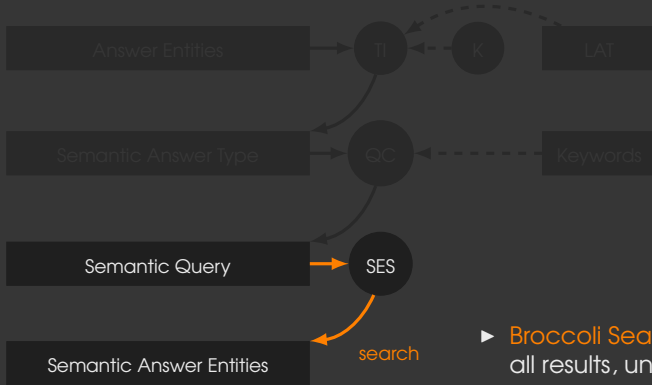
► Based on **type frequency**

# SEMANTIC QUERY CONSTRUCTION



- ▶ Rule-based:
  - \$1 is-a SAT
  - \$1 occurs-with keywords

# SEMANTIC SEARCH



► **Broccoli Search API:**  
all results, unfiltered



# ENTITY RETRIEVAL

## STRICT MATCHING RESULTS

	<b>Filter</b>		<b>R</b>	<b>P</b>	<b>R@S</b>	<b>P@S</b>	<b>F@S</b>
ont	qsim	ctype	(%)	(%)	(%)	(%)	(%)
			62	3	17	16	14
•			54	6	19	27	20
	•		61	3	36	32	30
		•	51	5	14	16	13
•	•		53	6	32	43	32
•		•	42	7	15	23	16
	•	•	56	6	35	38	32
•	•	•	48	8	30	44	31

Average results with **ontology filter** (ont), **query similarity** filter (qsim) and **coarse type filter** (ctype)

# ENTITY RETRIEVAL

## APPROXIMATE MATCHING RESULTS

	<b>Filter</b>		<b>R</b>	<b>P</b>	<b>R@S</b>	<b>P@S</b>	<b>F@S</b>
ont	qsim	ctype	(%)	(%)	(%)	(%)	(%)
			<b>88</b>	4	28	30	24
•			72	7	27	39	28
	•		87	4	<b>49</b>	49	43
		•	76	7	25	31	23
•	•		70	7	41	57	41
•		•	57	9	22	36	24
	•	•	78	7	48	59	<b>47</b>
•	•	•	63	<b>11</b>	39	<b>61</b>	41

Average results with **ontology filter** (ont), **query similarity** filter (qsim) and **coarse type filter** (ctype)

## SELECTION OPTIMALITY

Matching Type	$F@S_{opt}$ (%)	$R@S$ (%)	$P@S$ (%)	$F@S$ (%)	$Q_s$ (%)
strict	45	43	34	34	78
approximate	65	59	49	48	71

Selection quality compared to the **optimal selection**  $S_{opt}$

# TWO-PHASE APPROACH RESULTS

Phase	Matching Type	<i>R</i> (%)	<i>P</i> (%)	<i>P@R</i> (%)	<i>R@S</i> (%)	<i>P@S</i> (%)	<i>F@S</i> (%)
ER	strict	56	6	38	33	38	31
ER	approximate	78	7	56	47	60	46
DS	strict	44	9	24	20	22	19
DS	approximate	54	12	31	27	31	25

Overall results for both phases.

# CONCLUSION

- ▶ **Competitive results** in entity retrieval phase
  - ▶ Simple and effective filtering
  - ▶ Near-optimal selection method
  - ▶ High noise in entity extraction
- ▶ **Unsatisfactory** deep search results
  - ▶ Unreliable semantic type detection
  - ▶ Ignored relation between entities

# FUTURE WORK

- ▶ **Further optimize results** in entity retrieval phase
  - ▶ Add document segmentation
  - ▶ Increase number of retrieved documents
  - ▶ More robust named entity extraction
  - ▶ Enable entity linking
- ▶ **Improve** semantic query construction
  - ▶ Semantic type classification based on *Freebase*
  - ▶ Rule-based semantic type detection
    - ▶ "Who" → person
    - ▶ "Where" → location

# FUTURE WORK

- ▶ **Further optimize results** in entity retrieval phase
  - ▶ Add document segmentation
  - ▶ Increase number of retrieved documents
  - ▶ More robust named entity extraction
  - ▶ Enable entity linking
- ▶ **Improve** semantic query construction
  - ▶ Semantic type classification based on *Freebase*
  - ▶ Rule-based semantic type detection
    - ▶ "Who" → person
    - ▶ "Where" → location

Leverage existing **semantic search framework**

# PYTHIA

## SEMANTIC SEARCH ORACLE

### Quote

“For all the things we have to learn before we can do them,  
we learn by doing them.” *Aristotle*

### Repository

[github.com/eamsen/pythia](https://github.com/eamsen/pythia)



# ENTITY RANKING

## OVERALL

### Formula

$$\text{score}(e) = \sum_{s \in \text{Subscores}} \frac{w_s \cdot s(e)}{s_{\max}} \quad s_{\max} = \max_{n \in E} s(n)$$
$$\text{Subscores} = \{s_C, s_H, s_{CD}, s_{HD}\}$$

### Description

- ▶  $w_s$ : weighting parameter
- ▶  $s_C$ : document entity freq.
- ▶  $s_H$ : snippet entity freq.
- ▶  $s_{CD}$ : documents freq.
- ▶  $s_{HD}$ : snippets freq.

# ENTITY RANKING

## SUBSCORES

### Formula

$$s(e) = |Occurs(e)| \text{ for } s \in \{s_{CD}, s_{HD}\}$$

$$s(e) = \sum_{\langle freq, rank \rangle \in Occurs(e)} \frac{w_{rank} \cdot freq}{\log(cf(e) + cf_{base})} \text{ for } s \in \{s_C, s_H\}$$

### Description

- ▶  $w_{rank}$ : weighting constants
- ▶  $cf$ : corpus entity freq. (popularity)
- ▶  $cf_{base}$ : in range  $[1, \infty)$
- ▶  $\lambda$ : dampening parameter

$$w_{rank} = 1 - \frac{rank}{1 + \lambda \cdot rank_{max}}$$

# ANSWER SELECTION

## MOVING AVERAGE PIVOT

### Formula

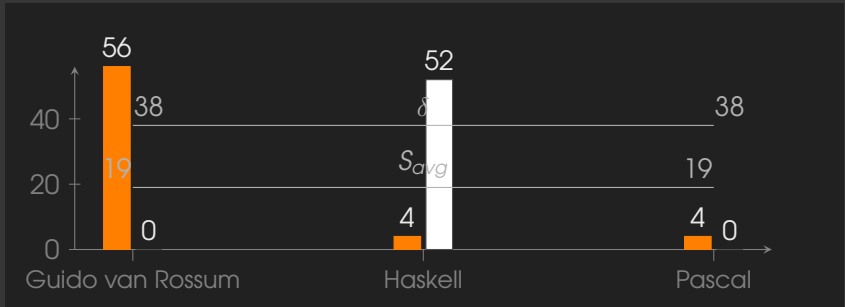
$$E_s = \{e \in E_c \mid e_s \geq \delta\} \quad \delta = S_{avg} + (2\gamma - 1)(S_{max} - S_{avg})$$

### Description

- ▶  $E_s$ : selection set
- ▶  $E_c$ : candidate set
- ▶  $e_s$ : entity score
- ▶  $\delta$ : score threshold
- ▶  $S_{avg_r} = \alpha \cdot e_{r-1_s} + (1 - \alpha) \cdot S_{avg_{r-1}}$  with  $S_{avg_1} = e_{1_s}$
- ▶  $\alpha = \frac{2}{|E_c|+1}$
- ▶  $\gamma$ : in range  $[0, 1]$

# ANSWER SELECTION

## EXAMPLE



Query "inventor of the python programming language":  
moving average score  $S_{avg} \approx 19$ , extrema  $S_{min} = 3$  and  
 $S_{max} = 83$ , with  $\gamma = 0.65$  we get the threshold  $\delta \approx 38$ .

# WHAT IS A NAMED ENTITY?

## Example

Milky Way, Mars, Alan Turing, *you*

## Properties

- ▶ Name
- ▶ Type
- ▶ Distinct identity

# TREC ENTITY TRACK: RELATED ENTITY FINDING

## Query

```
<query>
  <entity_name>Daft Punk</entity_name>
  <entity_url>daftpunk.com</entity_url>
  <target_entity>organisation</target_entity>
  <narrative>
    What recording companies sell Daft Punk songs?
  </narrative>
</query>
```

## Answer Records

virginrecords.com	1	0.98	.../wiki/Daft_Punk
somarecords.com	2	0.97	.../wiki/Daft_Punk
disney.go.com/music	3	0.89	.../wiki/Daft_Punk

# CONNECTION TO QUESTION ANSWERING

	Question Answering	Entity Retrieval
Emphasis	query semantics	entity ranking
Result type	factoid	entities
<i>in most cases</i>	factoids contain	entities