



清华大学  
Tsinghua University

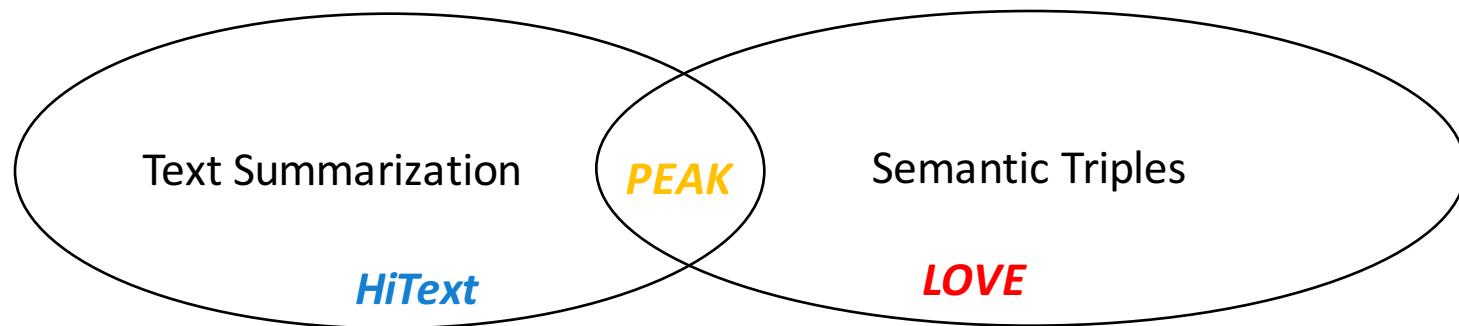
# Identifying and Connecting Salient Information via Semantic Text Representations

Qian Yang

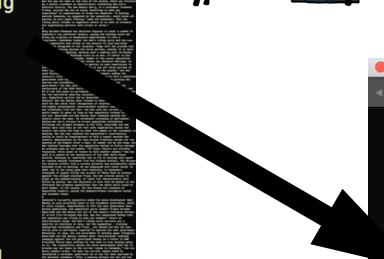
PhD Candidate, Tsinghua University

[www.larayang.com](http://www.larayang.com)

# Outline



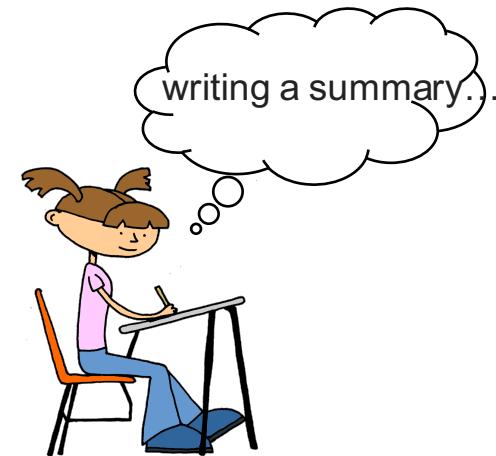
# Text Summarization



```
d30001.txt
```

1 Cambodian leader Hun Sen on Friday rejected opposition parties' demands for talks outside the country, accusing them of trying to "internationalize" the political crisis. Government and opposition parties have asked King Norodom Sihanouk to host a summit meeting after a series of post-election negotiations between the two opposition groups and Hun Sen's party to form a new government failed. Opposition leaders Prince Norodom Ranariddh and Sam Rainsy, citing Hun Sen's threats to arrest opposition figures after two alleged attempts on his life, said they could not negotiate freely in Cambodia and called for talks at Sihanouk's residence in Beijing. Hun Sen, however, rejected that. "I would like to make it clear that all meetings related to Cambodian affairs must be conducted in the Kingdom of Cambodia," Hun Sen told reporters after a Cabinet meeting on Friday. "No-one should internationalize Cambodian affairs. It is detrimental to the sovereignty of Cambodia," he said. Hun Sen's Cambodian People's Party won 64 of the 122 parliamentary seats in July's elections, short of the two-thirds majority needed to form a government on its own. Ranariddh and Sam Rainsy have charged that Hun Sen's victory in the elections was achieved through widespread fraud. They have demanded a thorough investigation into their election complaints as a precondition for their cooperation in getting the national assembly moving and a new government formed. Hun Sen said on Friday that the opposition concerns over their safety in the country was "just an excuse for them to stay abroad." Both

Line 1, Column 1 Tab Size: 4 Plain Text



```
d30001.m.100.t.a.sum
```

1 Prospects were dim for resolution of the political crisis in Cambodia in October 1998.  
2 Prime Minister Hun Sen insisted that talks take place in Cambodia while opposition leaders Ranariddh and Sam Rainsy, fearing arrest at home, wanted them abroad.  
3 King Sihanouk declined to chair talks in either place.  
4 A U.S. House resolution criticized Hun Sen's regime while the opposition tried to cut off his access to loans.  
5 But in November the King announced a coalition government with Hun Sen heading the executive and Ranariddh leading the parliament.  
6 Left out, Sam Rainsy sought the King's assurance of Hun Sen's promise of safety and freedom for all politicians.  
7

Line 1, Column 1 Tab Size: 4 Plain Text

- Automatically creating a compressed version of a given text that provides useful information for the user.

# Evaluation of Text Summarization

Evaluate summaries...

```
d30001.txt
```

1 Cambodian leader Hun Sen on Friday rejected opposition parties' demands for talks outside the country, accusing them of trying to "internationalize" the political crisis. Government and opposition parties have asked King Norodom Sihanouk to host a summit meeting after a series of post-election negotiations between the two opposition groups and Hun Sen's party to form a new government failed. Opposition leaders Prince Norodom Ranariddh and Sam Rainsy, citing Hun Sen's threats to arrest opposition figures after two alleged attempts on his life, said they could not negotiate freely in Cambodia and called for talks at Sihanouk's residence in Beijing. Hun Sen, however, rejected that. "I would like to make it clear that all meetings related to Cambodian affairs must be conducted in the Kingdom of Cambodia," Hun Sen told reporters after a Cabinet meeting on Friday. "No-one should internationalize Cambodian affairs. It is detrimental to the sovereignty of Cambodia," he said. Hun Sen's Cambodian People's Party won 64 of the 122 parliamentary seats in July's elections, short of the two-thirds majority needed to form a government on its own. Ranariddh and Sam Rainsy have charged that Hun Sen's victory in the elections was achieved through widespread fraud. They have demanded a thorough investigation into their election complaints as a precondition for their cooperation in getting the national assembly moving and a new government formed. Hun Sen said on Friday that the opposition concerns over their safety in the country was "just an excuse for them to stay abroad." Both

Line 1, Column 1      Tab Size: 4      Plain Text

homework



```
d30001.m.100.t.a.sum
```

1 Prospects were dim for resolution of the political crisis in Cambodia in October 1998.  
2 Prime Minister Hun Sen insisted that talks take place in Cambodia  
3 Cambodian prime minister Hun Sen rejects demands of 2 opposition  
4 parties for talks in Beijing after failing to win a 2/3 majority  
5  
6 Cambodia King Norodom Sihanouk praised formation of a coalition of  
7 the Countries top two political parties, leaving strongman Hun Sen  
8 as Prime Minister and opposition leader Prince Norodom Ranariddh  
9 president of the National Assembly.  
10  
11 The announcement comes after months of bitter argument following  
12 the failure of any party to attain the required quota to form a  
13 government.  
14  
15 Opposition leader Sam Rainey was seeking assurances that he and  
16 his party members would not be arrested if they return to Cambodia.  
17  
18 Rainey had been accused by Hun Sen of being behind an  
19 assassination attempt against him during massive street  
20 demonstrations in September.

Line 1, Column 1      Tab Size: 4      Plain Text

The picture of a teacher is from <http://www.clipartpanda.com/categories/teacher-clip-art>.

# Evaluating Summary Content

- **Human assessors**
  - Judge each summary individually
  - Very **time-consuming** and does not **scale** well
- **ROUGE** (Lin 2004)
  - Automatically compares **n-grams** from **target summaries** to **reference summaries (model summaries)** and can be applied to a large number of summaries.
  - Recall-based measure
    - ROUGE-1: 
$$\frac{\text{count of unigrams in reference that appear in system}}{\text{count of unigrams in reference summary}}$$

# N-grams

Physician note **“...Patient has evidence of macular degeneration...”**

Unigrams      “**patient**” “**has**” “**evidence**” “**of**” “**macular**” “**degeneration**”

Bigrams      “**patient has**” “**evidence of**” “**macular degeneration**”  
                “**has evidence**” “**of macular**”

Trigrams      “**patient has evidence**” “**of macular degeneration**”  
                “**has evidence of**”  
                “**evidence of macular**”

4-grams      “**patient has evidence of**”  
                “**has evidence of macular**”  
                “**evidence of macular degeneration**”

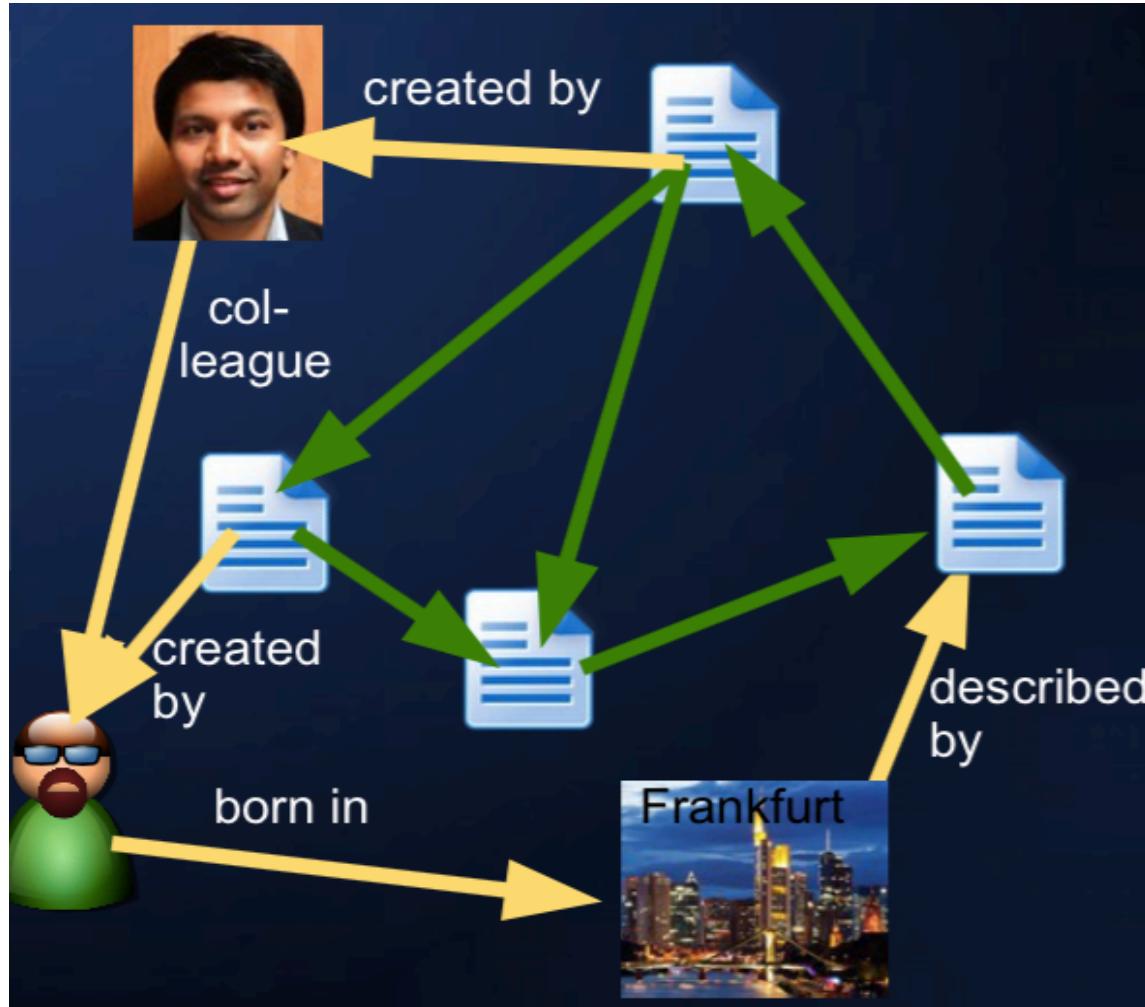
# Problem

- **N-grams** can be quite different between the reference summary and the target summary
  - "Benz developed a new kind of motor"
  - vs. "This new engine type was designed by Karl Benz"
- Our idea: Use **semantic triples** to evaluate summary content automatically

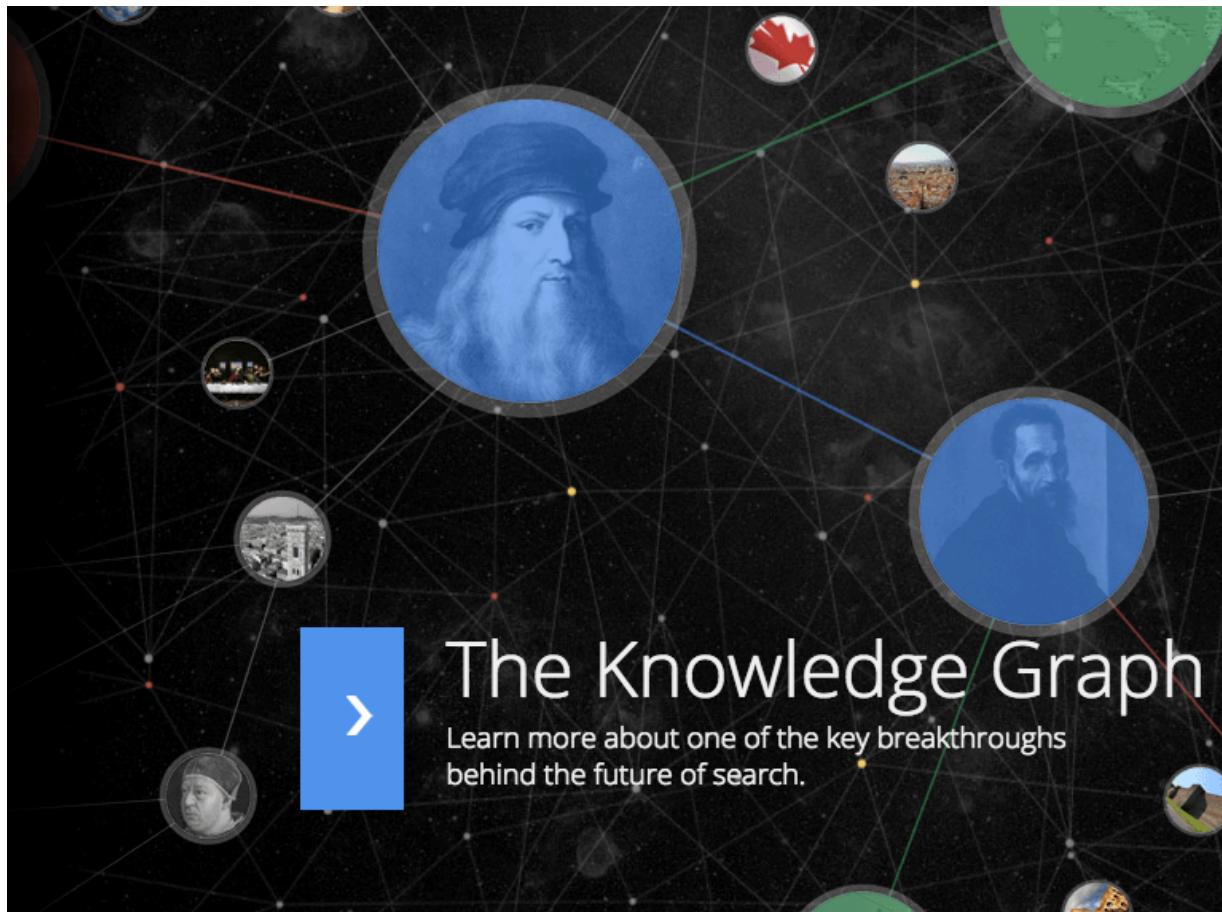
# Semantic Triples

- A semantic triple consists of (head entity, relation, tail entity)
  - Example: (City:Beijing, CityCapitalOfCountry, Country:China)

# Example: Semantic Web



# Example: Google Knowledge Graph



# Example: Google Knowledge Graph

Google Da Vinci

All Images Books News Videos More ▾ Search tools

About 80,100,000 results (0.46 seconds)

**Leonardo da Vinci - Wikipedia, the free encyclopedia**  
[https://en.wikipedia.org/wiki/Leonardo\\_da\\_Vinci](https://en.wikipedia.org/wiki/Leonardo_da_Vinci) ▾ Wikipedia ▾  
Leonardo di ser Piero da Vinci, more commonly Leonardo da Vinci or simply Leonardo was an Italian polymath whose areas of interest included invention, ...  
Mona Lisa - Vitruvian Man - The Last Supper - Personal life

**da Vinci Surgery - Minimally Invasive Robotic Surgery with ...**  
[www.davincisurgery.com/](http://www.davincisurgery.com/) ▾  
da Vinci Surgery - Minimally Invasive Robotic Surgery with the da Vinci Surgical System.

**Leonardo da Vinci - Artist, Mathematician, Inventor, Writer ...**  
[www.biography.com/people/leonardo-da-vinci-40396](http://www.biography.com/people/leonardo-da-vinci-40396) ▾  
Leonardo da Vinci was a leading artist and intellectual of the Italian Renaissance who's known for his enduring works "The Last Supper" and "Mona Lisa." ... Leonardo da Vinci - Mini Biography (TV-PG; 3:14) Leonardo da Vinci began apprenticing under the artist Verrocchio.

**Blackmagic Design: DaVinci Resolve 12**  
[https://www.blackmagicdesign.com/.../davinciresolve](http://www.blackmagicdesign.com/.../davinciresolve) ▾ Blackmagic Design ▾  
DaVinci Resolve 12.5 combines professional non-linear video editing with the world's most advanced color corrector so now you can edit, color correct, finish ...

**Leonardo da Vinci - Facts & Summary - HISTORY.com**  
[www.history.com/topics/leonardo-da-vinci](http://www.history.com/topics/leonardo-da-vinci) ▾ History ▾



**More images**

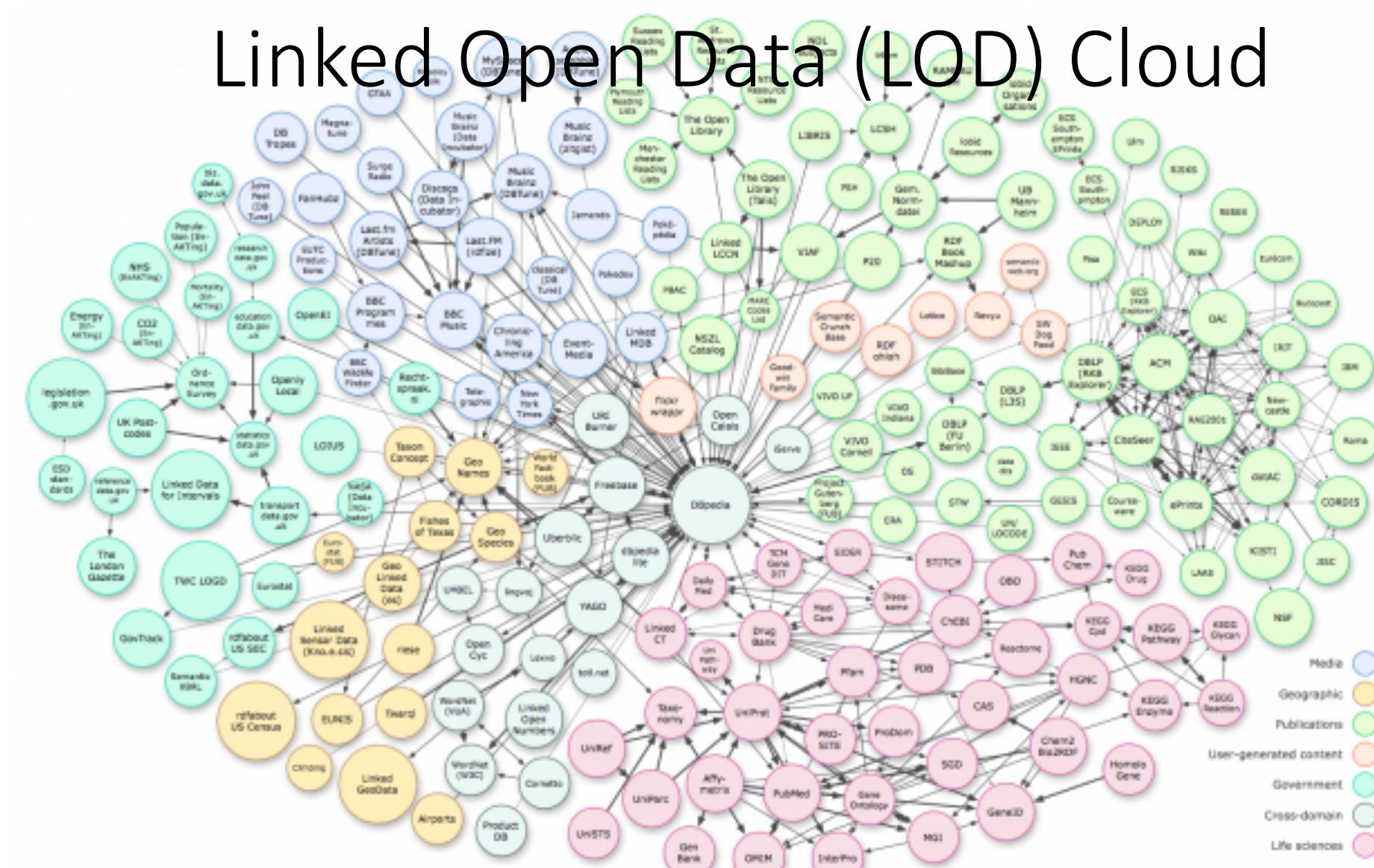
**Leonardo da Vinci** Share

Engineer

Leonardo di ser Piero da Vinci, more commonly Leonardo da Vinci or simply Leonardo, was an Italian polymath whose areas of interest included invention, painting, sculpting, architecture, science, music, ...  
[Wikipedia](#)

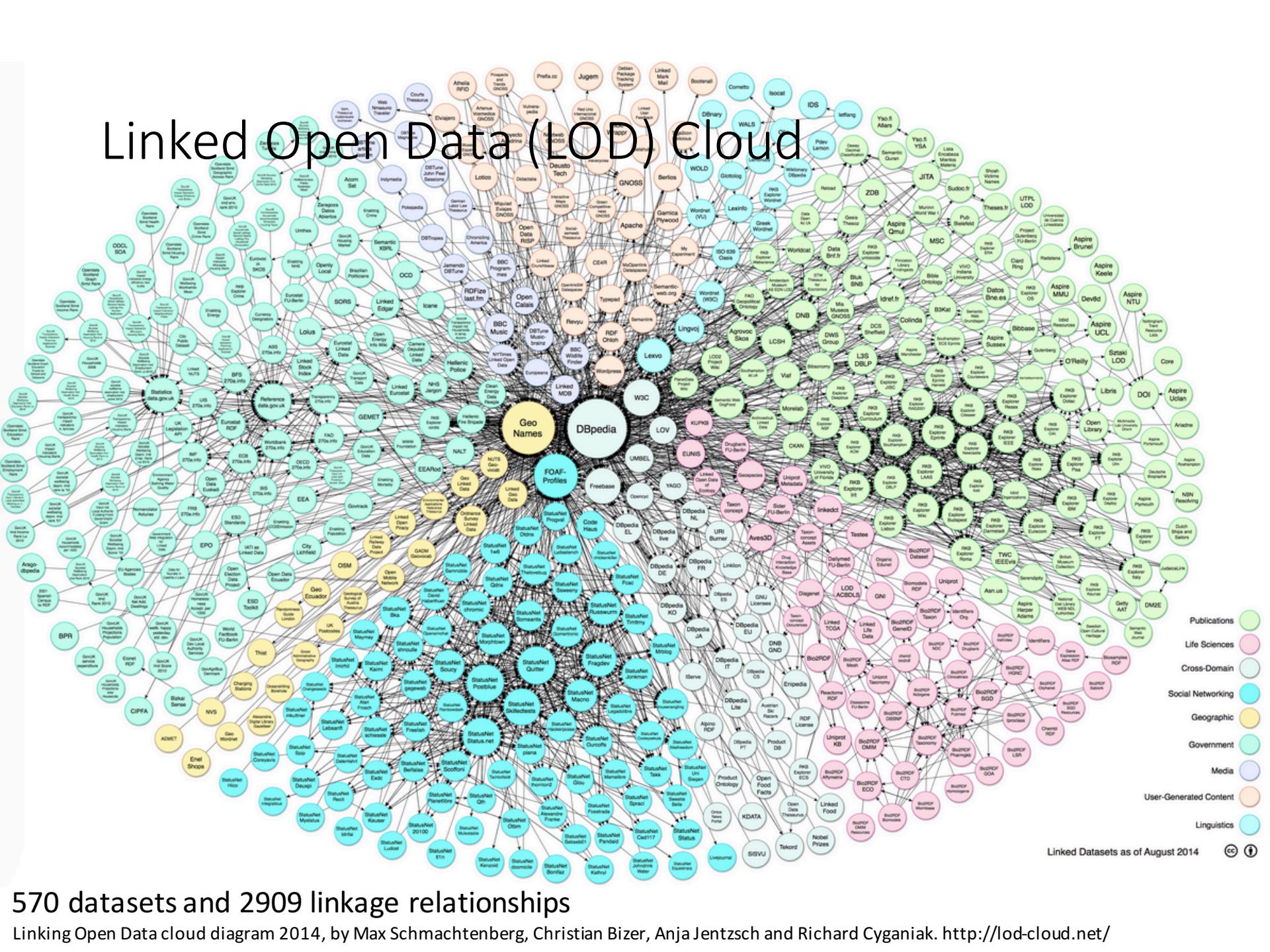
**Born:** April 15, 1452, Vinci, Italy  
**Died:** May 2, 1519, Amboise, France  
**Influenced:** Raphael, Peter Paul Rubens, Bernardino Luini, Antonio da Correggio, Pontormo  
**Siblings:** Bartolomeo da Vinci, Giovanni Ser Piero, More

# Linked Open Data (LOD) Cloud



As of September 2010

# Linked Open Data (LOD) Cloud



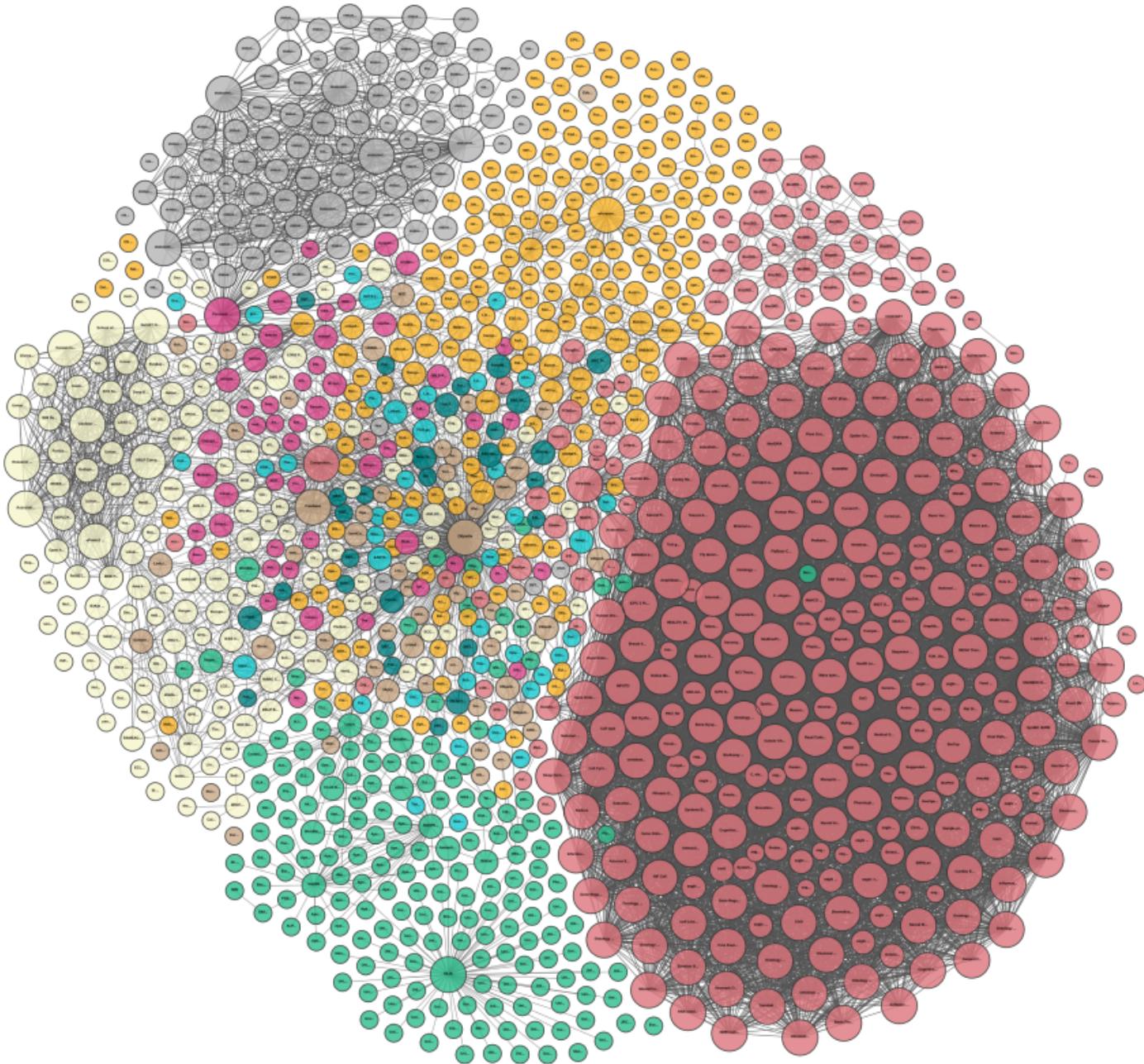
Linked Datasets as of August 2014



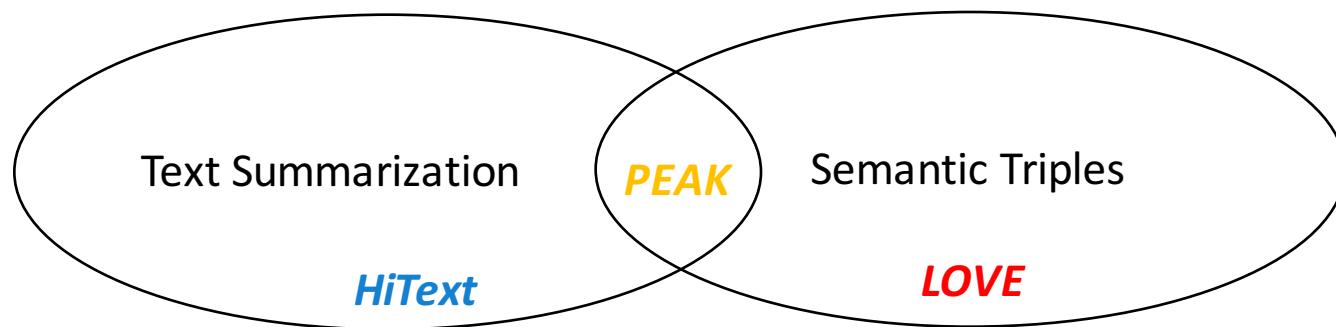
570 datasets and 2909 linkage relationships

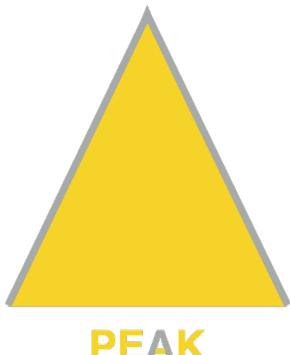
Linking Open Data cloud diagram 2014, by Max Schmachtenberg, Christian Bizer, Anja Jentzsch and Richard Cyganiak. <http://lod-cloud.net/>

Legend
Cross Domain
Geography
Government
Life Sciences
Linguistics
Media
Publications
Social Networking
User Generated
Incoming Links
Outgoing Links



# Outline





# PEAK: Pyramid Evaluation via Automated Knowledge Extraction

## 1. PEAK: Pyramid Evaluation via Automated Knowledge Extraction.

Qian Yang, Rebecca J. Passonneau, Gerard de Melo.

In: *Proc. AAAI 2016*. AAAI Press.

## 2. Wise Crowd Content Assessment and Educational Rubrics.

Rebecca J. Passonneau, Ananya Poddar, Gaurav Gite, Alisa Krivokapic, Qian Yang, Dolores Perin.

*International Journal of Artificial Intelligence in Education*, 10.1007/s40593-016-0128-6, 2016.

# Related Work

- **ROUGE** (Lin 2004)
  - Cannot capture similarity of meaning when little lexical overlap
- **Pyramid Method** (Nenkova and Passonneau, 2004)
  - Semantic comparison, reliable for individual summaries
  - Has required **manual** annotation
- **Our Method**
  - **No need** for **manually** created pyramids
  - Also good results on automatic assessment given a manually created pyramid

# Overview

- Semantic Content Analysis
- Pyramid Induction
- Automated Scoring of Summaries

# Semantic Content Analysis

## Model Summaries

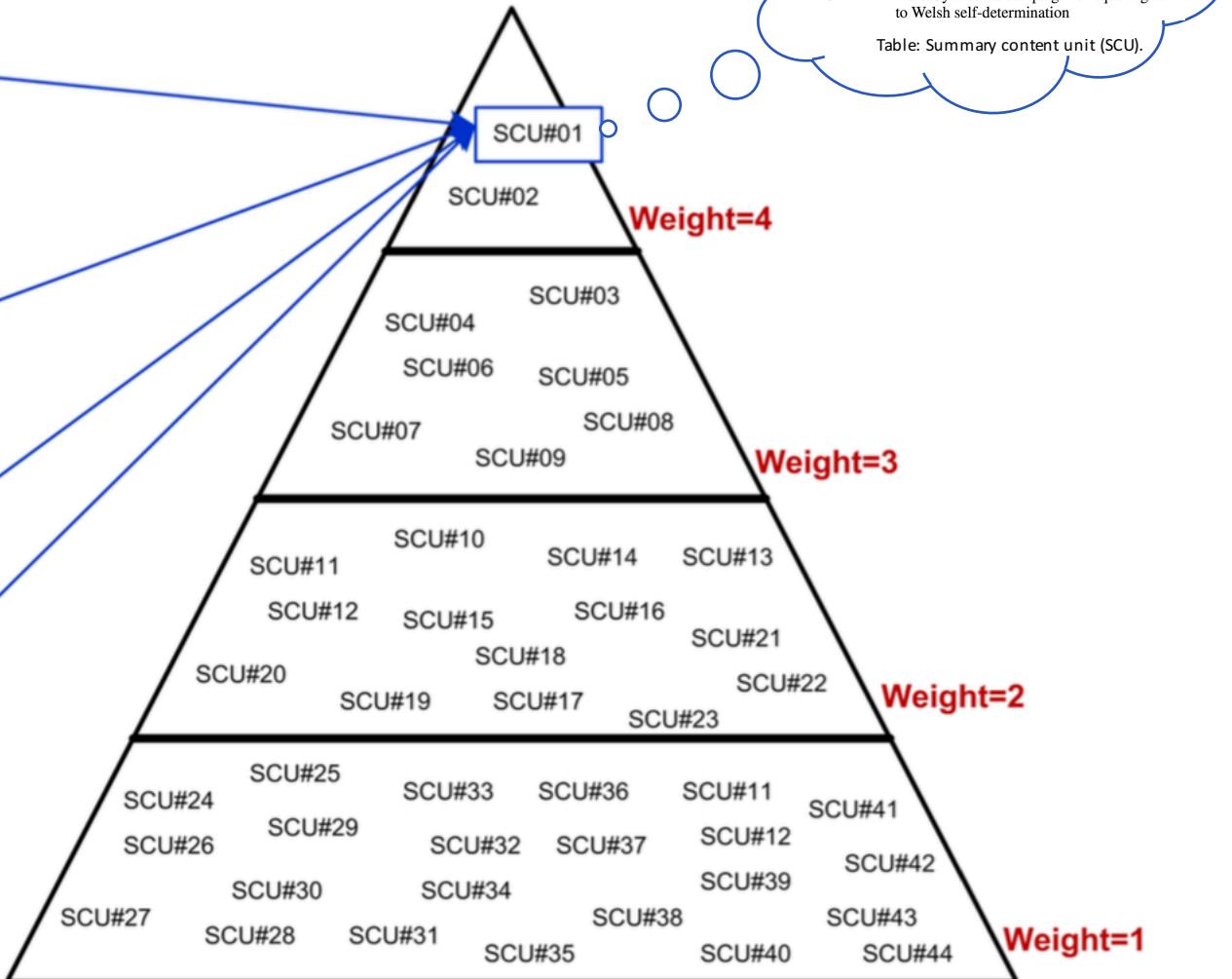
**Matter is what makes up all objects and substances**, and contains both volume and mass. Some types of matter are easily observable ...

The author of this passage titled 'What is Matter?' defines **matter as 'the stuff' that all objects and substances in the universe are made of.** ...

The passage, What is Matter, mainly focused on the topic of matter and its components. **Matter is identified as being present everywhere and in all substances.** ..

**Matter is all the objects and substances** that take up space **around us**. Matter can be detected and measured because it ...

## Pyramid Content Model



# Semantic Content Analysis

SCU	Plaid Cymru wants full independence
C1	Plaid Cymru wants full independence
C2	Plaid Cymru...whose policy is to...go for an independent Wales within the EC
C3	calls by...(Plaid Cymru)...fully self-governing Wales within the EC
C4	Plaid Cymru...its campaign for equal rights to Welsh self-determination

Weight: 4

Figure 1: Sample SCU from *Pyramid Annotation Guide: DUC 2006*.

# Semantic Content Analysis -PEAK

- “*The law of conservation of energy is the notion that energy can be transferred between objects but cannot be created or destroyed.*”
- Open information extraction (Open IE) methods split them and extract <subject,predicate,object> triples
- Example:

“*These characteristics determine the properties of matter*”  
yields the triple  
*⟨These characteristics, determine, the properties of matter⟩*
- We use ClausIE (Del Corro and Gemulla 2013)

# Semantic Content Analysis-PEAK

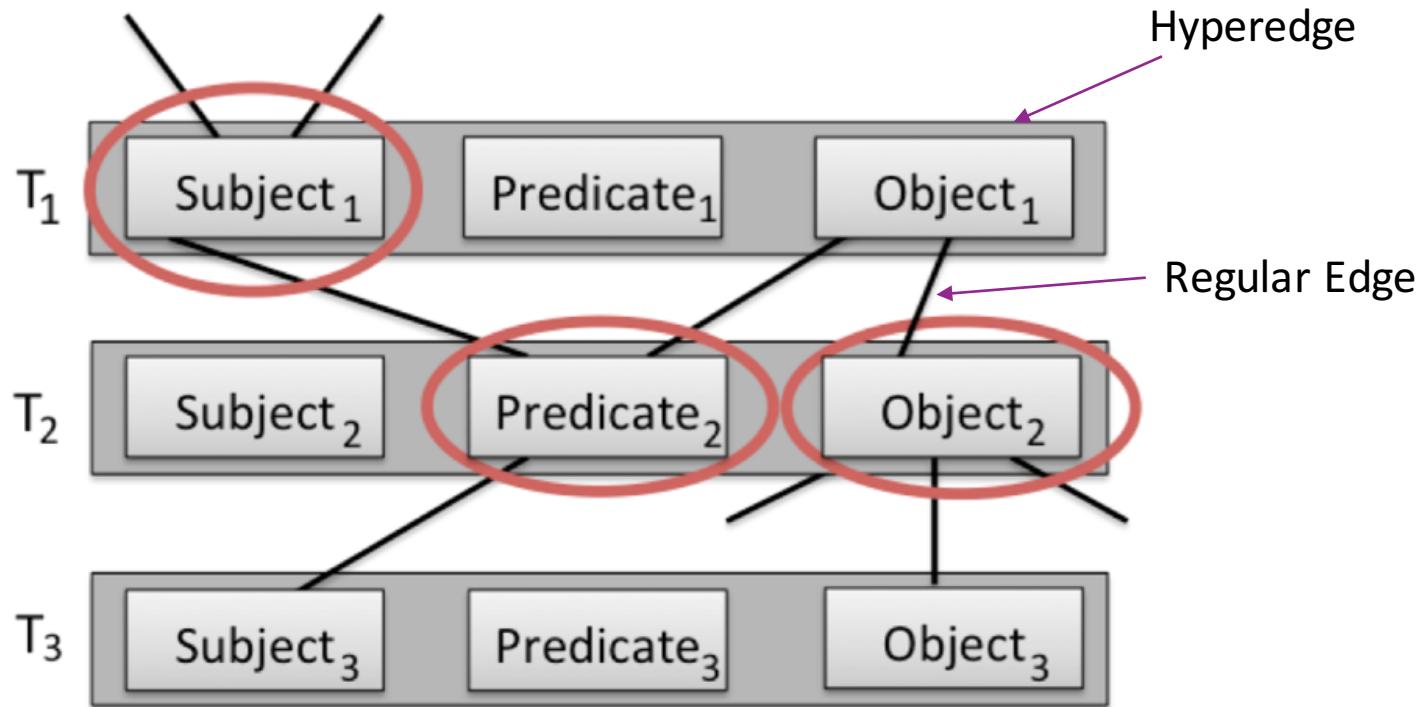


Figure 2: Hypergraph to capture similarities between elements of triples, with salient nodes circled in red

**Similarity Score:** Align, Disambiguate and Walk (ADW) (Pilehvar, Jurgens, and Navigli 2013)

# Pyramid Induction

**ANCHOR:** “**Matter**” “is” “all the objects and substances”

Similarity Class E1 for “Matter”	Similarity Class E2 for “all the objects and substances”
<p>“Matter”, “All matter”, “matter”, “the matter itself”, “a different matter”, “the matter itself systematically”, ...</p>	<p>“all the objects and substances”, “the substance”, “all objects and substances in the universe”, “what makes up all objects or substances and contains both volume and mass”, “as being present everywhere and in all substances”, ...</p>

# Pyramid Induction

**WEIGHT:** 4

**ANCHOR:** “Matter” “is” “all the objects and substances”

**FROM SENTENCE1 of CONTRIBUTOR1:** Matter is all the objects and substances that take up space around us.

**SALIENT NODES:** “Matter” “all the objects and substances”

**CONTRIBUTOR1:** “Matter” “is” “all the objects and substances”

**CONTRIBUTOR2:** “Matter” “is identified” “as being present everywhere and in all substances”

**CONTRIBUTOR3:** “The author of this passage titled What is Matter” “defines” “matter as the stuff that all objects and substances in the universe are made of”

**CONTRIBUTOR4:** “Matter” “is” “what makes up all objects or substances and contains both volume and mass”

# Pyramid Induction

---

## Algorithm 1 Merge similar SCUs

---

```
1: procedure MERGE(SCU anchors, weights)
2:   set a graph  $G$  whose nodes are all SCU anchors
3:   set threshold  $T_1$ 
4:   for each node  $anchor_m$  do
5:     for each node  $anchor_n$  do
6:       calculate  $similarityScore_{m,n}$ 
7:       if  $similarityScore_{m,n} \geq T_1$  then
8:         add edge between  $anchor_m$  and  $anchor_n$ 
9:   mergedSCU  $\leftarrow$  the connected component in  $G$ 
10:  mergedWeight  $\leftarrow$  max. weight of connected component
11:  Return  $mergedAnchor, mergedWeight$ 
```

---

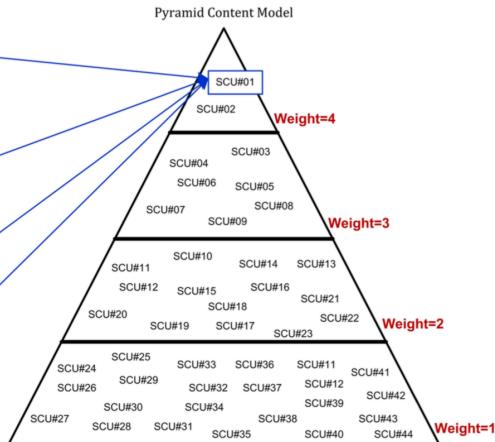
Model Summaries

Matter is what makes up all objects and substances, and contains both volume and mass. Some types of matter are easily observable ...

The author of this passage titled 'What is Matter?' defines matter as 'the stuff' that all objects and substances in the universe are made of ...

The passage, What is Matter? is concerned on the topic of matter and its components. Matter is identified as being present everywhere and in all substances...

Matter is all the objects and substances that take up space around us. Matter can be detected and measured because it ...



SCU 49

Plaid Cymru wants full independence

C1

Plaid Cymru wants full independence

C2

Plaid Cymru...whose policy is to...go for an independent Wales within the EC calls by...(Plaid Cymru)...fully

C3

self-governing Wales within the EC

C4

Plaid Cymru...its campaign for equal rights to Welsh self-determination

SCU created by the original Pyramid method



<b>WEIGHT: 4</b>
<b>ANCHOR:</b> "Matter" "is" "all the objects and substances"
<b>FROM SENTENCE1 of CONTRIBUTOR1:</b> Matter is all the objects and substances that take up space around us. <b>SALIENT NODES:</b> "Matter" "all the objects and substances"
<b>CONTRIBUTOR1:</b> "Matter" "is" "all the objects and substances" <b>CONTRIBUTOR2:</b> "Matter" "is identified" "as being present everywhere and in all substances" <b>CONTRIBUTOR3:</b> "The author of this passage titled What is Matter" "defines" "matter as the stuff that all objects and substances in the universe are made of" <b>CONTRIBUTOR4:</b> "Matter" "is" "what makes up all objects or substances and contains both volume and mass"

SCU created by PEAK

# Scoring – Pyramid Method

- Score a target summary against a pyramid
  - Annotators mark spans of text in the **target summary** that express an **SCU**
  - The SCU weights increment the raw score for the target summary
- Example:
  - **SCU Label:** Plaid Cymru wants full independence
  - **Target Summary:** **Plaid Cymru demands an independent Wales**

# Automated Scoring – PEAK

---

## Algorithm 2 Computing scores for target summaries

---

```
1: procedure SCORE(target summary  $sum$ )
2:   for each sentence  $s$  in  $sum$  do
3:      $T_s \leftarrow$  triples extracted from  $s$ 
4:     for each triple  $t \in \bigcup T_s$  do
5:       for each SCU  $s$  with weight  $w$  do
6:          $m \leftarrow$  similarity score between  $t$  and  $s$ 
7:         if  $m \geq T$  then
8:            $W[t][s] \leftarrow w$                                  $\triangleright$  store weight
9:    $S \leftarrow$  Munkres-Kuhn (Hungarian) Algorithm( $W$ )
10:  Return  $S$ 
```

---

- Score a target summary against a pyramid
  - Annotators mark **spans** of text in the **target summary** that express an **SCU**
  - The SCU weights increment the raw score for the target summary




---

### Algorithm 2 Computing scores for target summaries

---

```

1: procedure SCORE(target summary sum)
2:   for each sentence s in sum do
3:      $T_s \leftarrow$  triples extracted from s
4:     for each triple t  $\in \bigcup T_s$  do
5:       for each SCU s with weight w do
6:          $m \leftarrow$  similarity score between t and s
7:         if  $m \geq T$  then
8:            $W[t][s] \leftarrow w$                                  $\triangleright$  store weight
9:    $S \leftarrow$  Munkres-Kuhn (Hungarian) Algorithm(W)
10:  Return S

```

---

# Results: Student Summaries

- Dataset:
  - Twenty target summaries written by students from Perin et al. (2013).
  - Five reference model summaries from Passonneau et al. (2013).
- Metric: Pearson's correlation
- P: PEAK's pyramid v.s. P1 and P2: two human pyramids
- A: Our automatic method v.s. M: Manual method

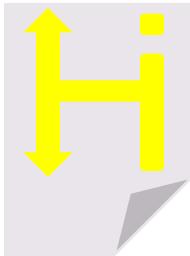
	P1 + M. Scoring	P2 + M. Scoring
P + A.Scoring	0.8263	0.7769
P2 +A. Scoring	0.8538	0.8112
P1 + M. Scoring	1	0.8857

# Results

- Machine-Generated Summaries
  - Dataset: the 2006 Document Understanding Conference (DUC) administered by NIST (“DUC06”)
  - The Pearson’s correlation score between PEAK’s scores and the manual ones is 0.7094.

# Conclusion

- The first fully automatic version of the pyramid method
- Not only evaluates target summaries but also generates the pyramids automatically
- Experiments show that
  - Our SCUs are similar to those created by humans
  - The method for assessing target summaries automatically has a high correlation with human assessors



# HiText: Text Reading with Dynamic Salience Marking

**HiText: Text Reading with Dynamic Salience Marking**

Qian Yang, Yong Cheng, Sen Wang, Gerard de Melo.

In: *Proc. WWW 2017 (Digital Learning Track)*. ACM.

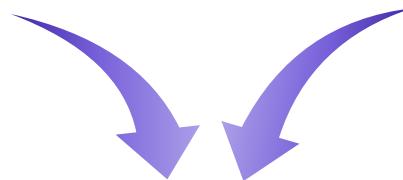
# Content

- Chanllenge and Related Work
- The HiText Method
- Interface Evaluation
- General Discussion and Conclusion

# Challenge



- Information Overload
  - Massive, non-stop 24/7 stream of new information
- Human ability
  - People can't peruse all information



- Current solution
  - Macro-level content selection
    - Aggregators (Google News), Recommendation engines
    - Helpful but insufficient



The figure above is from <http://chiefexecutive.net/ceos-can-reduce-information-overload/>  
The figure below is from google news.

# Our Focus

- Micro-level selection *within* a given document to help people reading text more efficient
  - Typical: hit-or-miss form of skimming
    - Glancing at somewhat **randomly** selected parts of the text
    - More efficient but not perfect
  - Our HiText Approach: Help the readers by specially marking on the screen those parts of the text that are likely to be **salient**.

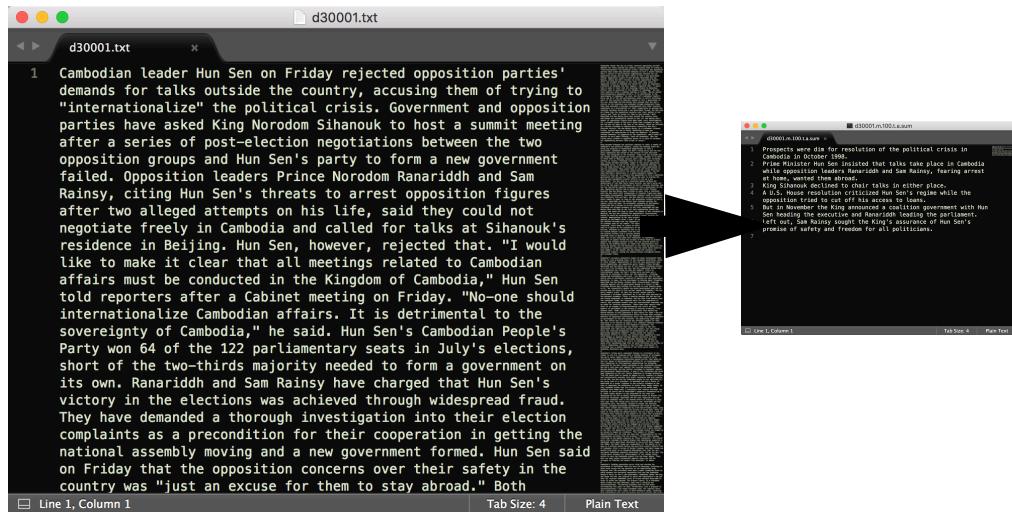
# Related Work

- Short Excerpts/Snippets
  - Problem: missing context.

A screenshot of a Google search results page for "Tsinghua University". The results include:

- Tsinghua University**  
www.tsinghua.edu.cn/publish/newthue/ ▾  
Gu Yuxiu: Tsinghua, where the Heart and Soul Belongs 2017.3.31 ... After entering Tsing Hua School (the predecessor of Tsinghua University) in 1915, he tried ... Schools & Departments · International Students · Graduate · Admissions
- Tsinghua University**  
www.tsinghua.edu.cn/ ▾ Translate this page  
电话查号台: 010-62793001 管理员信箱: webmaster@tsinghua.edu.cn 地址: 北京市海淀区清华大学 京公网安备110402430053号. 版权所有© 清华大学访问 ...  
You've visited this page 2 times. Last visit: 1/13/17
- (IIIS), Tsinghua University - 清华大学**  
www.iiis.tsinghua.edu.cn/en/ ▾  
Yingmei Liu 2017.05.19; Samuel Marshall Ganzfried 2017.04.29; Martin Charles Golumbic 2017.04.20; Bing Qi 2017.04.20; P.C. Ching 2017.03.24; Tan Lee  
You've visited this page 2 times. Last visit: 12/7/16
- Tsinghua University - Wikipedia**  
https://en.wikipedia.org/wiki/Tsinghua\_University ▾  
Tsinghua University is a research university located in Beijing, China, established in 1911. With strong research and training, Tsinghua University is consistently ... History · Academics · Student life · Campus

- Summaries
  - Problem: Too many details missing, hard to see the rest of document



# Related Work

- Receiving a summary is unsatisfactory
  - Automatic summarization systems are known to make mistakes [Gillick 2011]
  - Different readers may exhibit different interests [Rapp 2005]
  - Certain elements are vital for comprehension, such as
    - text structure [Carrell 1985]
    - formatting
    - figures and tables [Yi 2014]

# Related Work

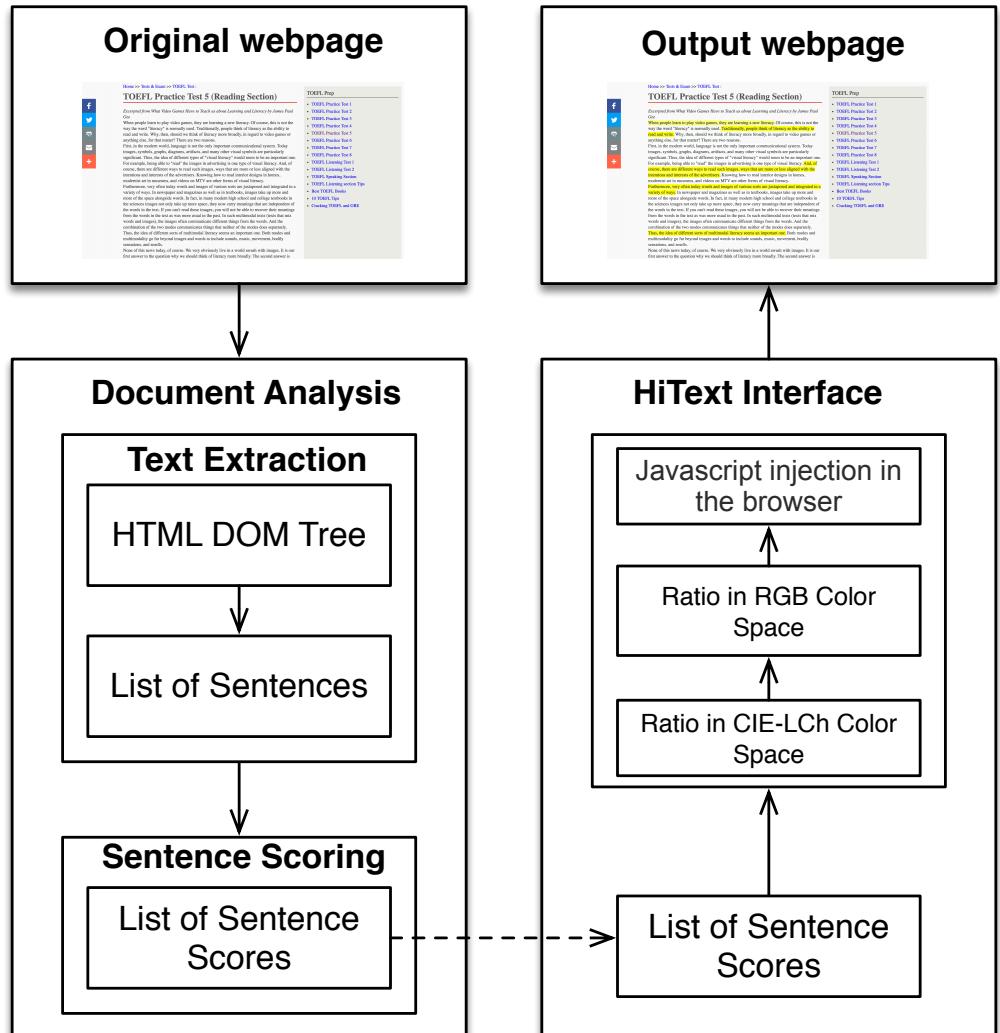
- Reading and comprehension
  - Three forms of behavior that readers seem to combine when skimming under time pressure
    - (1) scanning
    - (2) satisficing, i.e., skipping ahead, possibly to the next paragraph, once the information gain drops below a threshold
    - (3) sampling
  - Salience-based marking could lead to gains in efficiency

# Related Work

- The approach we follow is to highlight the key sentences in a text.
- Reading Assistance
  - The Semantize system [Wecker 2014]
  - A study in [Paiva 2014]
  - The ScentHighlights system [Chi 2005]
  - The Nestor Highlighter Extension

# The HiText Method

- 1. Document Analysis
  - Text Extraction
  - Sentence Scoring
- 2. HiText Interface
  - Document Rendering
  - Top phrase highlighting
  - Dynamic graded highlighting



System Overview

# 1. Document Analysis

- **Text Extraction**

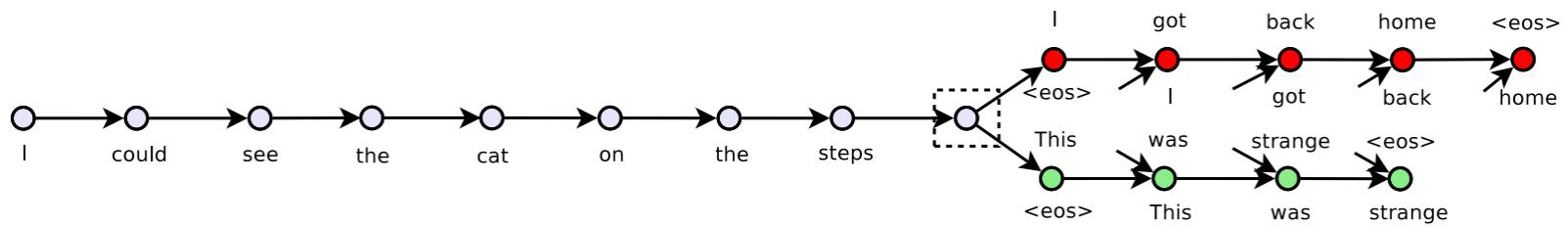
- Parse the input HTML, get HTML DOM tree
- Identify the main text using a stop word ratio-based heuristic
- Split each paragraph
  - Ignore abbreviations, such as U.K.

- **Sentence Scoring**

- Each extracted sentence  $s$  is analysed and assessed using a **deep learning-based scoring technique** to produce a salience score  $\sigma(s)$

# Skip-Thought Vector Approach [Kiros, 2015]

Given a tuple  $(s_{i-1}, s_i, s_{i+1})$  of contiguous sentences, with  $s_i$  the i-th sentence of a book, the sentence  $s_i$  is encoded and tries to reconstruct the previous sentence  $s_{i-1}$  and next sentence  $s_{i+1}$ .



In this example, the input is the sentence triplet *I got back home. I could see the cat on the steps. This was strange.* Unattached arrows are connected to the encoder output. Colors indicate which components share parameters. <eos> is the end of sentence token.

# Our Method

- We use a pre-trained model for this architecture to produce a 4,800-dimensional real-valued vector representation  $v_s$  of every sentence  $s$  in the document
- Compute initial sentence scores
  - $\sigma_0(s) = \sum_{s'} \frac{v_s^t v_{s'}'}{\|v_s\| \|v_{s'}\|}$
- Final scores
  - $\sigma(s) = \max\{0, 1 - 2 \frac{r(\sigma_0(s)) - 1}{n}\}$

## 2. The HiText Interface

- Document rendering
  - Inject additional JavaScript code
- Top phrase highlighting
  - Highlights the top- $k$  sentences
- Dynamic graded highlighting
  - Capture the location of the pointing device
  - $\sigma_{\min} = \min_s \sigma(s)$ ,  $\sigma_{\max} = \max_s \sigma(s)$
  - Assign two background colors  $C_{\max}$ ,  $C_{\min}$  to  $\sigma_{\max}$ ,  $\sigma_{\min}$
  - $\frac{\sigma(s) - \sigma_{\min}}{\sigma_{\max} - \sigma_{\min}} (C_{\max} - C_{\min}) + C_{\min}$

# TOEFL Practice Test 5 (Reading Section)

*Excerpted from What Video Games Have to Teach us about Learning and Literacy* by James Paul Gee

When people learn to play video games, they are learning a new literacy. Of course, this is not the way the word "literacy" is normally used. Traditionally, people think of literacy as the ability to read and write. Why, then, should we think of literacy more broadly, in regard to video games or anything else, for that matter? There are two reasons.

First, in the modern world, language is not the only important communicational system. Today images, symbols, graphs, diagrams, artifacts, and many other visual symbols are particularly significant. Thus, the idea of different types of "visual literacy" would seem to be an important one. For example, being able to "read" the images in advertising is one type of visual literacy. And, of course, there are different ways to read such images, ways that are more or less aligned with the intentions and interests of the advertisers. Knowing how to read interior designs in homes, modernist art in museums, and videos on MTV are other forms of visual literacy.

Furthermore, very often today words and images of various sorts are juxtaposed and integrated in a variety of ways. In newspaper and magazines as well as in textbooks, images take up more and more of the space alongside words. In fact, in many modern high school and college textbooks in the sciences images not only take up more space, they now carry meanings that are independent of the words in the text. If you can't read these images, you will not be able to recover their meanings from the words in the text as was more usual in the past. In such multimodal texts (texts that mix words and images), the images often communicate different things from the words. And the combination of the two modes communicates things that neither of the modes does separately. Thus, the idea of different sorts of multimodal literacy seems an important one. Both modes and multimodality go far beyond images and words to include sounds, music, movement, bodily sensations, and smells.

# TOEFL Practice Test 5 (Reading Section)

*Excerpted from What Video Games Have to Teach us about Learning and Literacy* by James Paul Gee

When people learn to play video games, they are learning a new literacy. Of course, this is not the way the word "literacy" is normally used. Traditionally, people think of literacy as the ability to read and write. Why, then, should we think of literacy more broadly, in regard to video games or anything else, for that matter? There are two reasons.

First, in the modern world, language is not the only important communicational system. Today images, symbols, graphs, diagrams, artifacts, and many other visual symbols are particularly significant. Thus, the idea of different types of "visual literacy" would seem to be an important one. For example, being able to "read" the images in advertising is one type of visual literacy. And, of course, there are different ways to read such images, ways that are more or less aligned with the intentions and interests of the advertisers. Knowing how to read interior designs in homes, modernist art in museums, and videos on MTV are other forms of visual literacy.

Furthermore, very often today words and images of various sorts are juxtaposed and integrated in a variety of ways. In newspaper and magazines as well as in textbooks, images take up more and more of the space alongside words. In fact, in many modern high school and college textbooks in the sciences images not only take up more space, they now carry meanings that are independent of the words in the text. If you can't read these images, you will not be able to recover their meanings from the words in the text as was more usual in the past. In such multimodal texts (texts that mix words and images), the images often communicate different things from the words. And the combination of the two modes communicates things that neither of the modes does separately.

Thus, the idea of different sorts of multimodal literacy seems an important one. Both modes and multimodality go far beyond images and words to include sounds, music, movement, bodily sensations, and smells.

# TOEFL Practice Test 5 (Reading Section)

*Excerpted from What Video Games Have to Teach us about Learning and Literacy by James Paul Gee*

When people learn to play video games, they are learning a new literacy. Of course, this is not the way the word "literacy" is normally used. Traditionally, people think of literacy as the ability to read and write. Why, then, should we think of literacy more broadly, in regard to video games or anything else, for that matter? There are two reasons.

First, in the modern world, language is not the only important communicational system. Today images, symbols, graphs, diagrams, artifacts, and many other visual symbols are particularly significant. Thus, the idea of different types of "visual literacy" would seem to be an important one. For example, being able to "read" the images in advertising is one type of visual literacy. And, of course, there are different ways to read such images, ways that are more or less aligned with the intentions and interests of the advertisers. Knowing how to read interior designs in homes, modernist art in museums, and videos on MTV are other forms of visual literacy.

Furthermore, very often today words and images of various sorts are juxtaposed and integrated in a variety of ways. In newspaper and magazines as well as in textbooks, images take up more and more of the space alongside words. In fact, in many modern high school and college textbooks in the sciences images not only take up more space, they now carry meanings that are independent of the words in the text. If you can't read these images, you will not be able to recover their meanings from the words in the text as was more usual in the past. In such multimodal texts (texts that mix words and images), the images often communicate different things from the words. And the combination of the two modes communicates things that neither of the modes does separately.

Thus, the idea of different sorts of multimodal literacy seems an important one. Both modes and multimodality go far beyond images and words to include sounds, music, movement, bodily sensations, and smells.

# TOEFL Practice Test 5 (Reading Section)

*Excerpted from What Video Games Have to Teach us about Learning and Literacy by James Paul Gee*

When people learn to play video games, they are learning a new literacy. Of course, this is not the way the word "literacy" is normally used. Traditionally, people think of literacy as the ability to read and write. Why, then, should we think of literacy more broadly, in regard to video games or anything else, for that matter? There are two reasons.

First, in the modern world, language is not the only important communicational system. Today images, symbols, graphs, diagrams, artifacts, and many other visual symbols are particularly significant. Thus, the idea of different types of "visual literacy" would seem to be an important one. For example, being able to "read" the images in advertising is one type of visual literacy. And, of course, there are different ways to read such images, ways that are more or less aligned with the intentions and interests of the advertisers. Knowing how to read interior designs in homes, modernist art in museums, and videos on MTV are other forms of visual literacy.

Furthermore, very often today words and images of various sorts are juxtaposed and integrated in a variety of ways. In newspaper and magazines as well as in textbooks, images take up more and more of the space alongside words. In fact, in many modern high school and college textbooks in the sciences images not only take up more space, they now carry meanings that are independent of the words in the text. If you can't read these images, you will not be able to recover their meanings from the words in the text as was more usual in the past. In such multimodal texts (texts that mix words and images), the images often communicate different things from the words. And the combination of the two modes communicates things that neither of the modes does separately.

Thus, the idea of different sorts of multimodal literacy seems an important one. Both modes and multimodality go far beyond images and words to include sounds, music, movement, bodily sensations, and smells.

# Interface Evaluation

- Experiment 1: Comprehension Survey
- Experiment 2: Comprehension Tests
  - Experiment 2-A: Response Time
  - Experiment 2-B: Time-Restricted Case
- Experiment 3: Salience Scoring

# Experiment 1: Comprehension Survey

- Evaluating comprehension is non-trivial
  - Impossible to repeat a given task with different methods under equal conditions
- One way is by finding two participants with comparable reading abilities.
- In our experiment: more robust by relying on two groups (with counterbalancing)

# Experiment 1: Comprehension Survey

- Participants
  - A pool of 10 participants
- Materials
  - Two articles from popular online journals
    - one 543 words, the other 1260 words
- Procedure and Measures
  - Randomly divided into two groups
  - Read a given text under two different conditions in a counterbalanced order
    - Top5, Top5+Graded

# Experiment 1: Comprehension Survey

	<b>Article 1</b>	<b>Article 2</b>
HiText Group	86 (40.37, 18.06)	187 (183.15, 81.91)
Control Group	127 (66.86, 29.90)	234 (121.57, 54.37)
Time saved	32.28%	20.09%

Table 1: Efficiency analysis (Experiment 1), given as means (std. deviation, std. error) in seconds.

	Strongly Agree	Agree	Not Sure	Disagree	Strongly Disagree
Q1	1.00	0.00	0.00	0.00	0.00
Q2	0.90	0.10	0.00	0.00	0.00
Q3	0.90	0.10	0.00	0.00	0.00

Table 2: Qualitative feedback from participants.

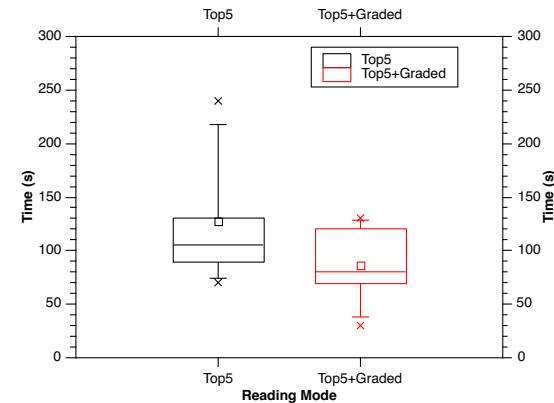
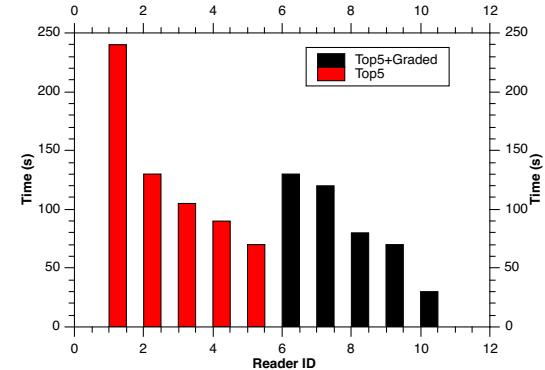


Figure 4: Analysis of reading times for article 1 (Experiment 1) as a bar chart and box plot. Top5 refers to the control group and serves as the baseline.

# Experiment 2: Comprehension Tests

- Experiment 2-A: Response Time
  - Materials (from TOEFL exams)
    - Text 1 (Plain Text) 649 words; Text 2 (Top5) 974 words; Text 3 (Top5+Graded) 818 words
    - One question to evaluate the reader's global understanding
  - Participants, Procedure, and Measure
    - Only participants that answered **all three questions correctly** were considered.
    - Until we had 10 participants with correct answers
    - Measure the **time**

# Experiment 2: Comprehension Tests

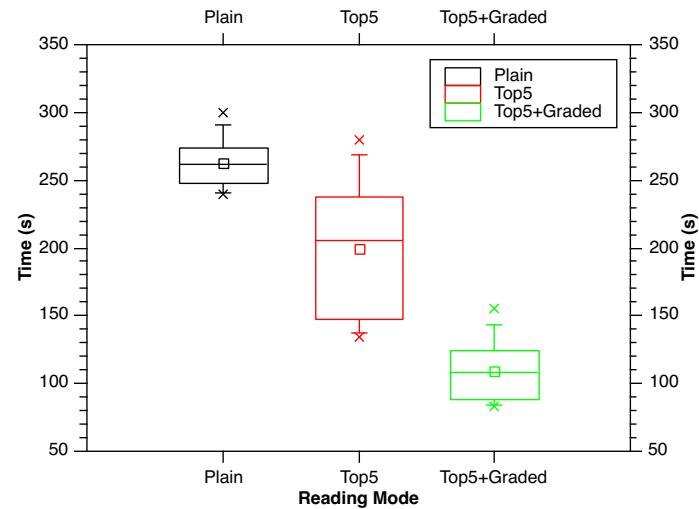
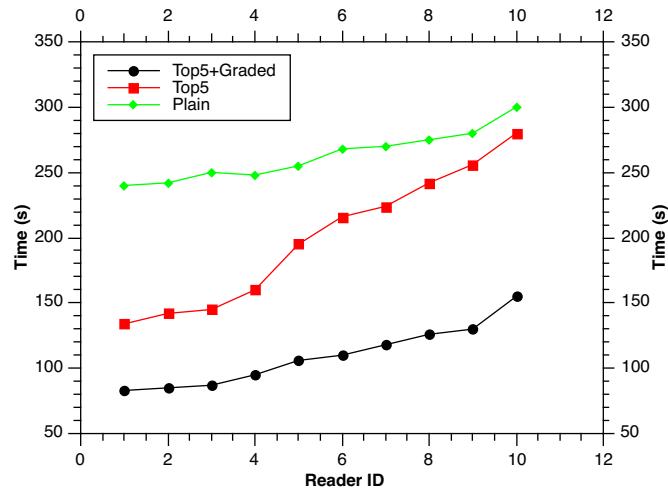


Figure 5: Response time for correctly answering reading comprehension questions (Experiment 2-A)

	Plain	Top5	Top5+Graded
Average Time (SD, SE)	262.8 (19.16, 6.06)	199.4 (52.17, 16.50)	109.5 (23.25, 7.35)
T-test (paired)	w/ Top5 $t = 5.928$ , p-value = 0.0002213	w/ Top5+Graded $t = 9.5539$ , p-value = 5.225e-06	w/ Plain $t = 89.832$ , p-value = 1.331e-14

Table 3: Efficiency analysis (Experiment 2-A), given as means (std. deviation, std. error) in seconds.

# Experiment 2: Comprehension Tests

- Experiment 2-B: Time-Restricted Case
  - Participants and Materials
    - A different set of 10 participants
    - The same materials as in Experiment 2-A.
  - Procedure
    - First read the texts and then answer the question **within 2 minutes**.
    - Respond to Question, **irrespective of** whether the given time is enough or not

# Experiment 2: Comprehension Tests

- Results and Discussion
- The resulting ratios of correct results for the three settings
  - For plain text, only 40% of participants answered the question correctly
  - For Top5 highlighting, 70% answered correctly
  - While for Top5 + Graded highlighting, in fact 100% made the right choice

# Experiment 3: Salience Scoring

- Materials
  - One of the texts from our pool of texts from Experiment 1.
- Participants and Procedure
  - Two independent annotators
    - The 7-point Likert-style importance scale given by Vagias [Vagias 2006]
  - For LexRank [Erkan 2004], SumBasic [Nenkova 2005], and Luhn [Luhn 1958]
    - Gradually increases the length of the expected summaries
    - Ultimately rank the salience of every sentence
  - For HiText
    - Obtain the rank directly from its salience scores.

# Experiment 3: Salience Scoring

	Annotator 1	Annotator 2
Annotator 2	0.91	1.00
Our method	0.75	0.83
LexRank	0.53	0.52
SumBasic	-0.06	-0.05
Luhn	-0.15	-0.06

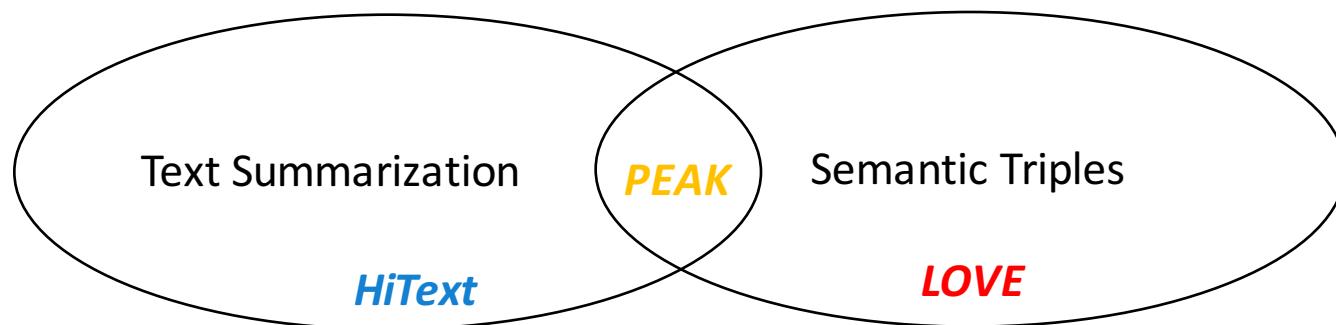
Table 4: Spearman correlations with human salience assessments.

- Deep learning-based representations
  - Learn to map different sentences to similar representations
- Word-based comparisons ( $\checkmark$  large document clusters,  $\times$  short text units)
  - Example: car and automobile

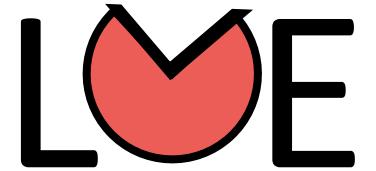
# CONCLUSION

- A simple but effective new method for guiding the reader towards  
**salient content units** in text
- Enables **faster and better** reading comprehension.

# Outline

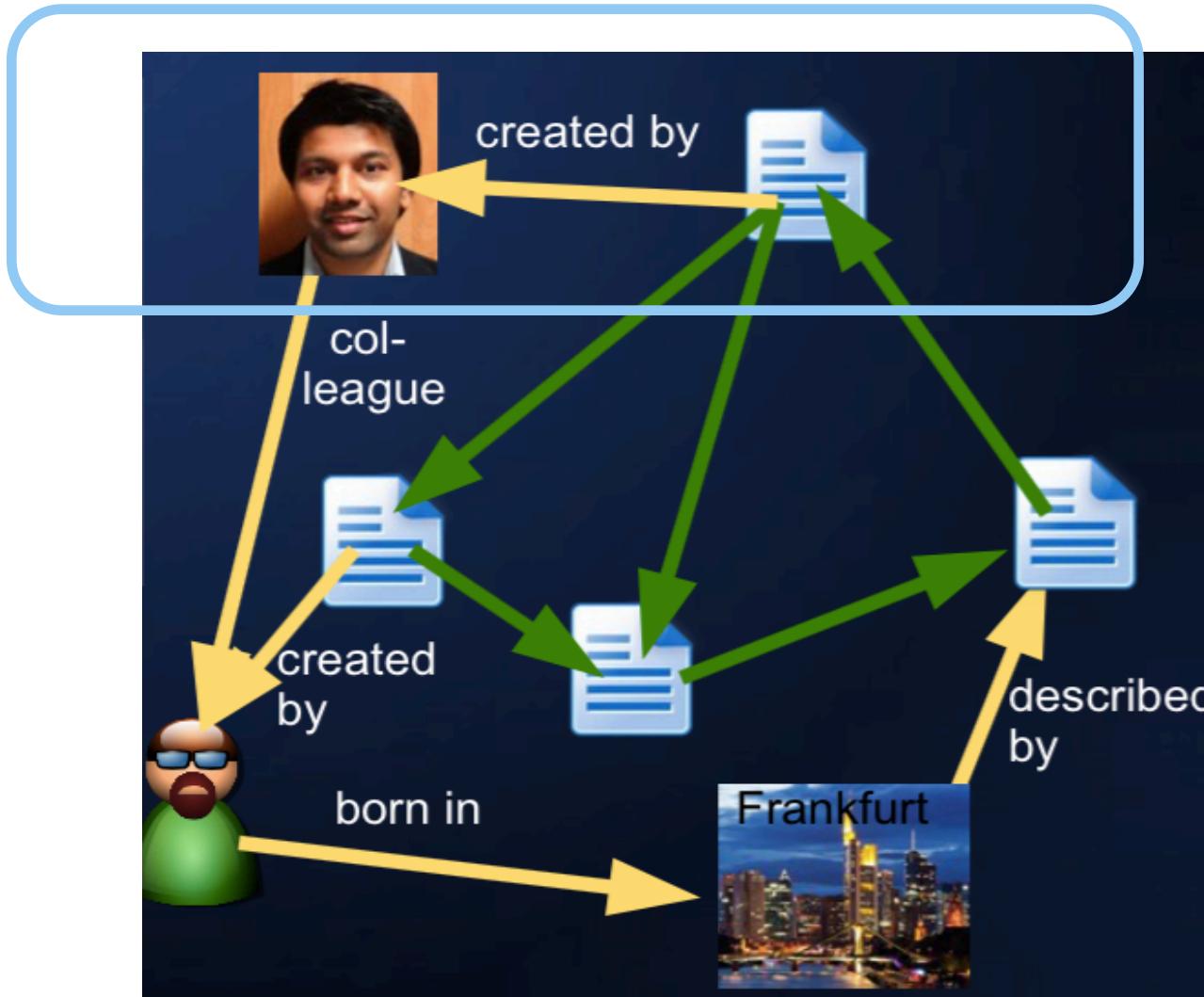


- Example
  - "The book's author is Kafka."  
"Kafka wrote the book."  
"Kafka created the book."



LOVE: Disambiguation and Alignment of Properties  
via Linked Open Vocabularies Enrichment

# Example: Semantic Web



# Problem



## Term Name: creator

URI:	<u><a href="http://purl.org/dc/elements/1.1/creator">http://purl.org/dc/elements/1.1/creator</a></u>
Label:	Creator
Definition:	An entity primarily responsible for making the resource.
Comment:	Examples of a Creator include a person, an organization, or a service. Typically, the name of a Creator should be used to indicate the entity.

# Problem

schema.org

Search

## Thing > Property > author

The author of this content. Please note that author is special in that HTML 5 provides a special mechanism for indicating authorship via the rel tag. That is equivalent to this and may be used interchangeably.

### Values expected to be one of these types

Organization

Person

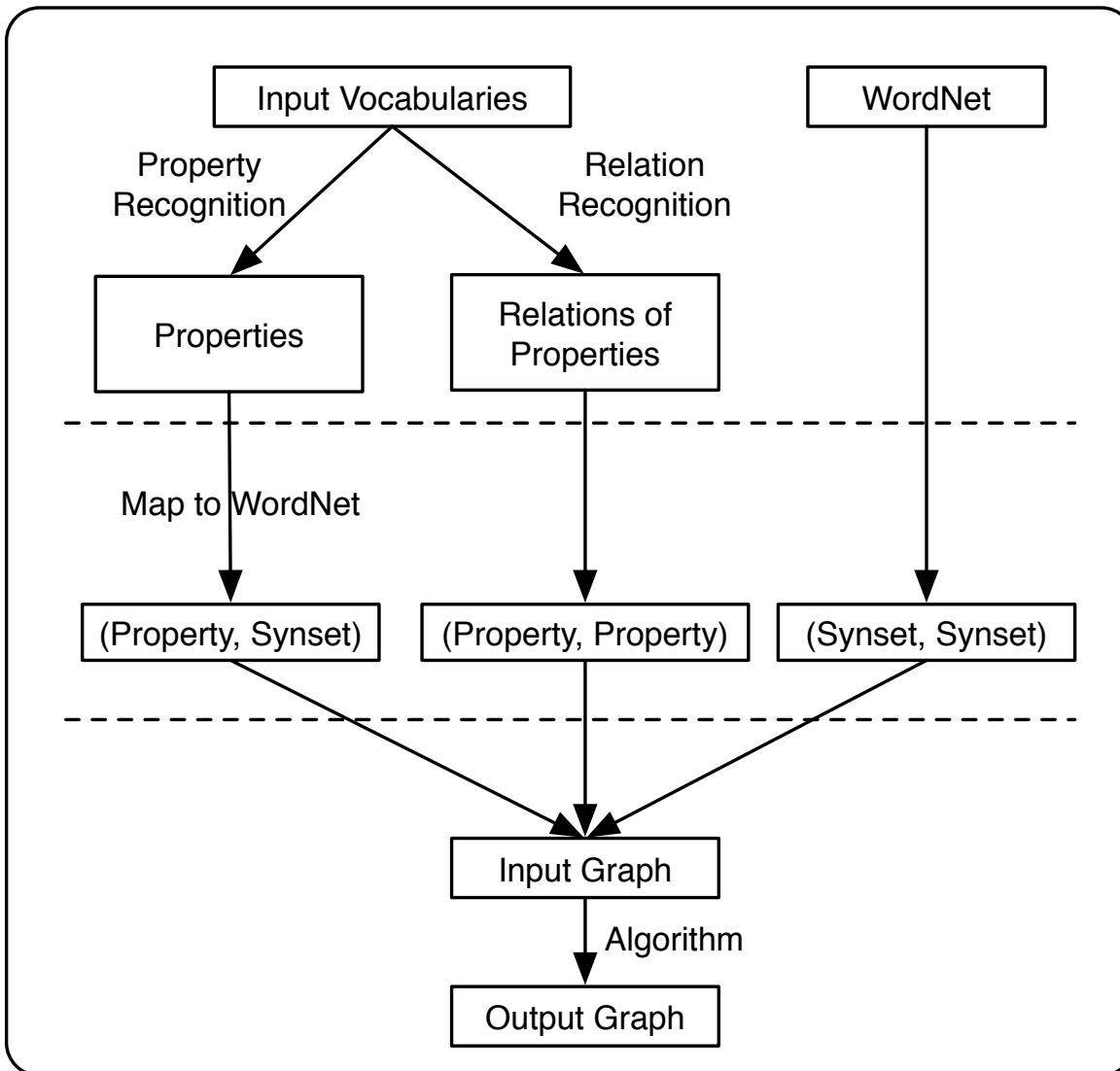
### Used on these types

CreativeWork

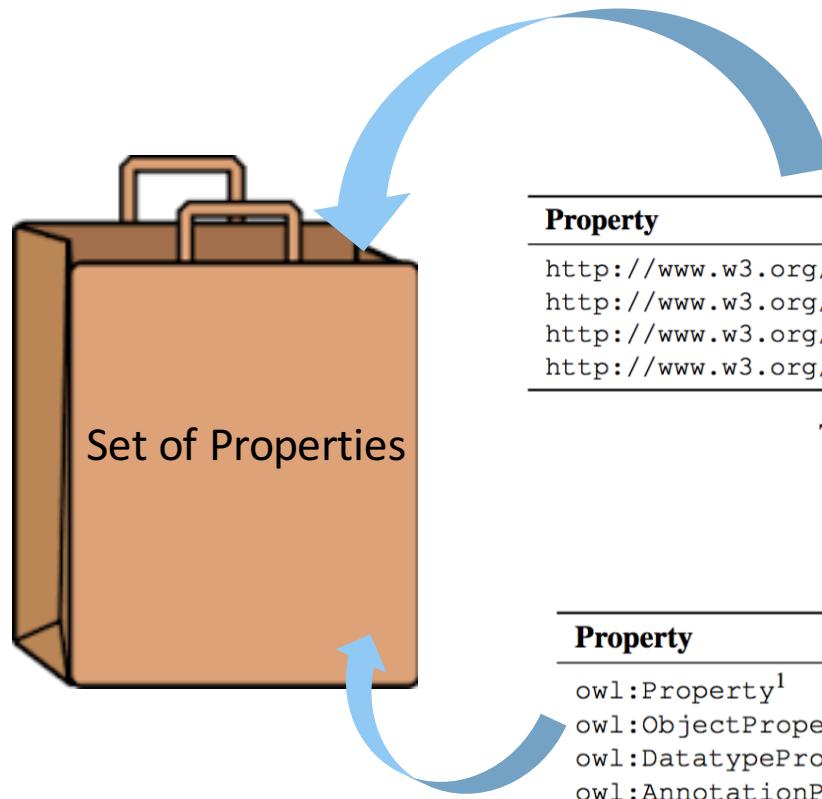
# Our goal is more ambitious

- PEAK: if triples are **similar**
- LOVE: if one is more **specific** than the other (subProperty) or **inverse** of the other

# Overview of our approach



# Collect all the properties



Property	Relation
<code>http://www.w3.org/2002/07/owl#equivalentProperty</code>	equivalent
<code>http://www.w3.org/2002/07/owl#propertyDisjointWith</code>	nonEquivalent
<code>http://www.w3.org/2002/07/owl#inverseOf</code>	inverse
<code>http://www.w3.org/2000/01/rdf-schema#subPropertyOf</code>	subProperty

Table 1: Predicates for four relations.

Property
<code>owl:Property</code> <sup>1</sup>
<code>owl:ObjectProperty</code>
<code>owl:DatatypeProperty</code>
<code>owl:AnnotationProperty</code>

Table 2: Types of Properties.

The picture of bag is from <http://cn.clipartlogo.com/free/brown-paper-bag-lunch.html>

---

**Algorithm 1** Collect all the properties

---

```
1: procedure COLLECT(input vocabularies in  $\mathcal{V}$ )
2:    $T \leftarrow \bigcup_{V \in \mathcal{V}} V$                                  $\triangleright$  input triples
3:    $P \leftarrow \{\}$                                           $\triangleright$  empty set
4:   for each triple  $t_i = \langle s_i, p_i, o_i \rangle \in T$  do
5:     if its predicate  $p_i$  is listed in Table 1 then
6:        $P \leftarrow P \cup \{s_i, o_i\}$ 
7:       if  $o_i$  is in Table 2 (with appropriate predicate) then
8:          $P \leftarrow P \cup \{s_i\}$ 
9:    $P \leftarrow P \cup$  items in Table 1
10:   $P \leftarrow P \cup$  items in Table 2
11:   $A \leftarrow \{t \in T \mid t \text{ describes attribute of } p \in P\}$   $\triangleright$  attributes
12:   $L \leftarrow \{t \in T \mid \text{property in Table 1}\}$             $\triangleright$  links
13:  return  $P, A, L$ 
```

---

# Mapping to Wordnet: Example



## Term Name: creator

URI:	<u><a href="http://purl.org/dc/elements/1.1/creator">http://purl.org/dc/elements/1.1/creator</a></u>
Label:	Creator
Definition:	An entity primarily responsible for making the resource.
Comment:	Examples of a Creator include a person, an organization, or a service. Typically, the name of a Creator should be used to indicate the entity.

[Synset: [Offset: 9637345] [POS: noun] Words: **creator** -- (a person who grows or makes or invents things)] [PointerType: [Label: hypernym] [Key: @] Applies To: noun, verb]]

[Synset: [Offset: 9559474] [POS: noun] Words: **Godhead, Lord, Creator, Maker, Divine, God Almighty, Almighty, Jehovah** -- (terms referring to the Judeo-Christian God)]

---

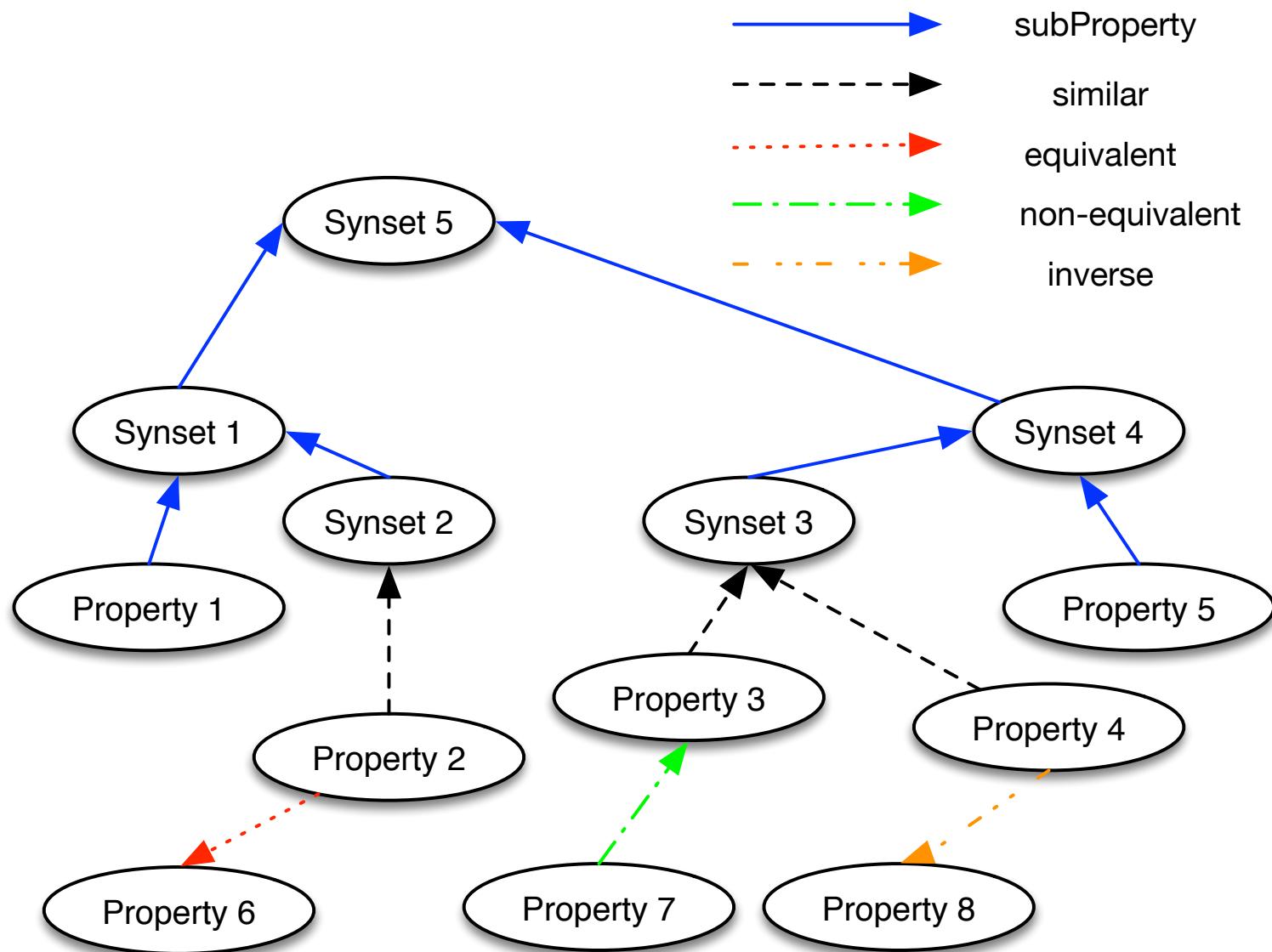
## Algorithm 2 Mapping to Wordnet

---

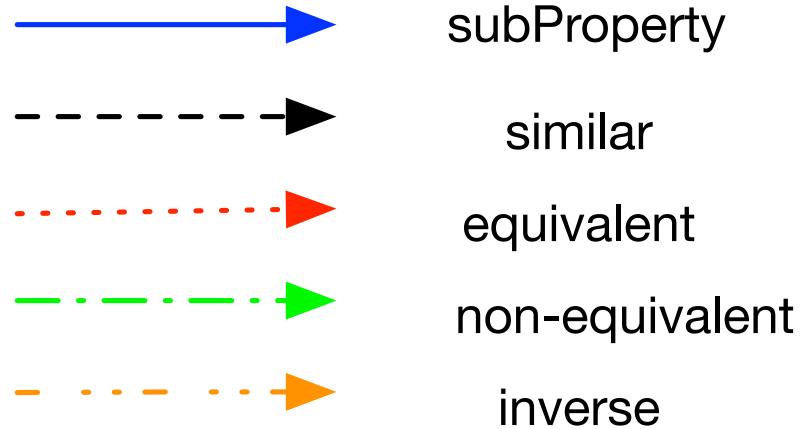
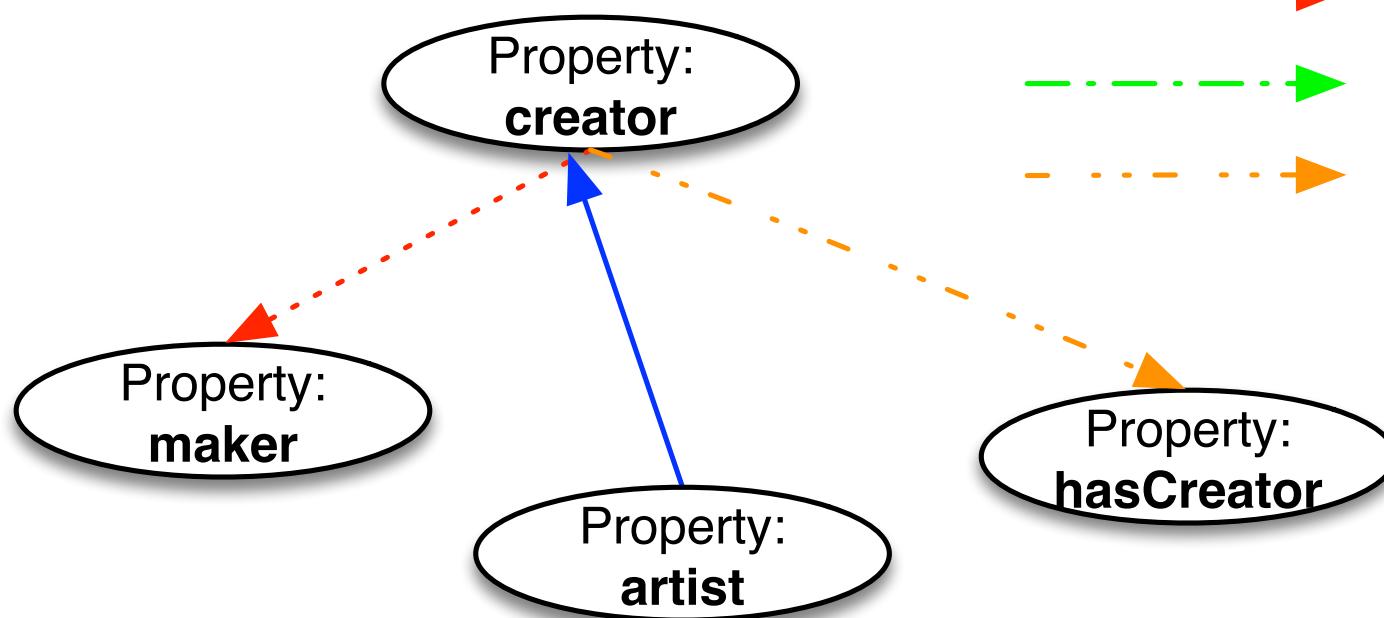
```
1: procedure MAP(set of properties  $P$ , WordNet synsets  $S$ )
2:    $V \leftarrow P \cup S$ ,  $E \leftarrow \emptyset$                                  $\triangleright$  Graph
3:   for each property  $p \in P$  do
4:     extract keyword  $k_p$  from  $p$                        $\triangleright$  Cf. definition in text
5:     obtain POS tag  $pos_p$  of  $k_p$                           $\triangleright$  Stanford CoreNLP
6:      $S_p \leftarrow$  synsets of  $k_p$  with  $pos_p$  in WordNet
7:     for every synset  $s_p \in S_p$  do
8:        $E \leftarrow E \cup \{\text{weighted edge between } p, s_p$ 
           $\text{with weight } score_{(c_p, s_p)}\}$ 
9:   return  $G = (V, E)$ 
```

---

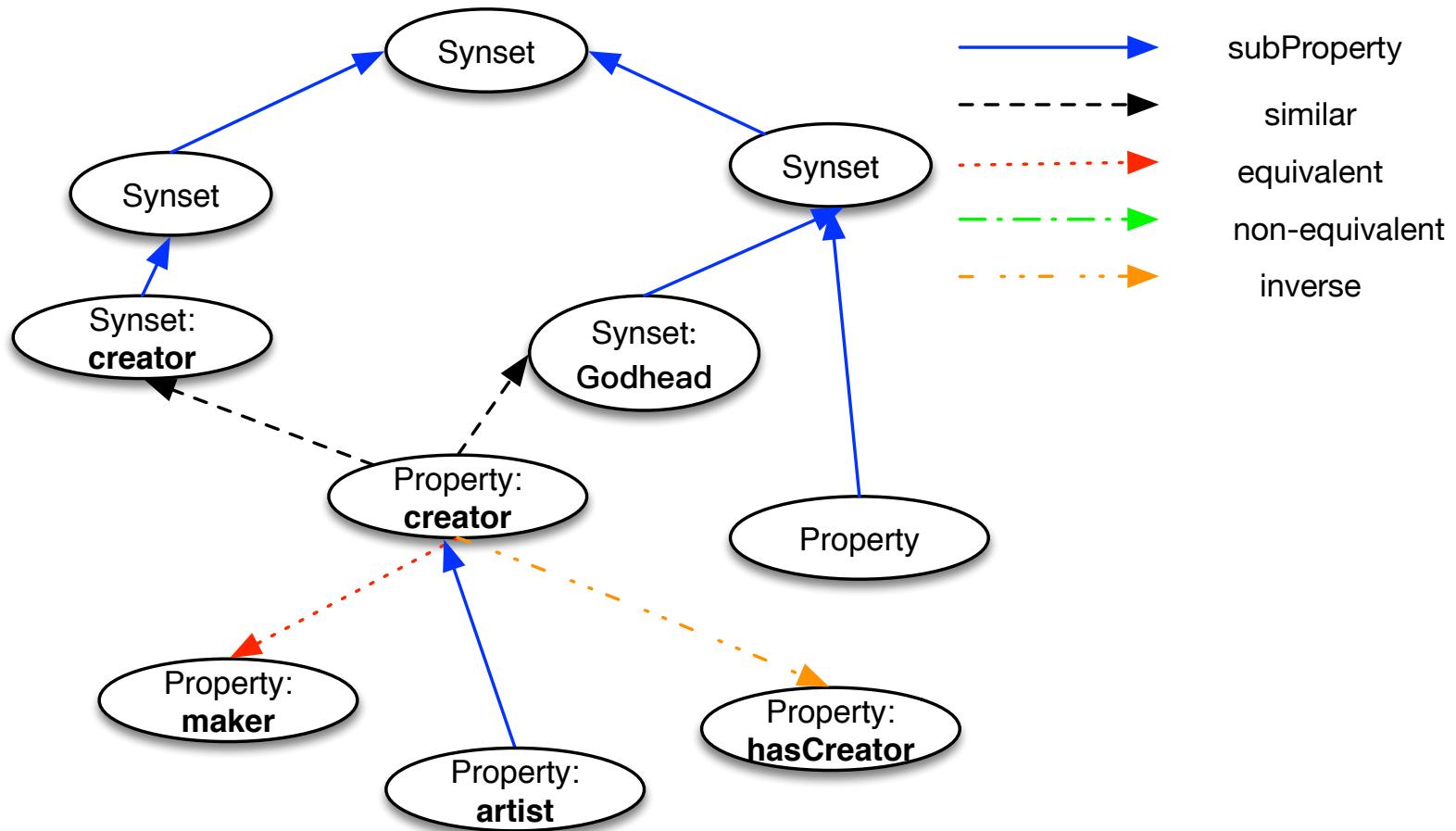
# Conceptual-graph of the Input Graph



# Example: Input Graph



# Example: Input Graph



# LP Problem

$$\underset{x}{\text{maximize}} \quad \sum_{p \in \Phi, m, n} x_{mn}^p + \sum_{p \in \phi, m} x_{mm}^p - \sum_{p \in \phi, m, n} \alpha_p * C * \epsilon_{mn}^p \quad (2a)$$

$$\text{subject to} \quad x_{mn}^p \leq x_{nm}^p \quad \forall m, n, p \in \{eq, frameEq, nonEq, inverseOf, hasE0P1, hasE0P3, hasE0P5\} \quad (2b)$$

$$x_{mn}^p \leq x_{mn}^q \quad \forall m, n, p \in \{eq, inverseOf, synFrameEq\}, q \in \{frameEq\} \quad (2c)$$

$$x_{mn}^p + x_{nl}^p - 1 \leq x_{ml}^p \quad \forall m, n, l, p \in \{eq, frameEq, subOf\} \quad (2d)$$

$$x_{mn}^p \leq 1 - x_{nm}^p \quad \forall m, n, p \in \{subOf\} \quad (2e)$$

$$x_{mn}^p + x_{nl}^p - 1 \leq x_{ml}^q + x_{ml}^r + x_{lm}^r \quad \forall m, n, l, p \in \{inverseOf\}, q \in \{eq\}, r \in \{subOf\} \quad (2f)$$

$$x_{mn}^p \leq 1 - x_{mn}^q \quad \forall m, n, p \in \{subOf\}, q \in \{eq, inverseOf\}$$

and  $\forall m, n, p \in \{nonEq\}, q \in \{eq\}$

and  $\forall m, n, p \in \{eq\}, q \in \{inverseOf\}$  (2g)

$$\forall m, n, l, p \in \{eq\}, q \in \{nonEq, subOf, inverseOf\} \quad (2h)$$

$$\forall m, n, p \in \{hasE0P1, hasE0P3, hasE0P5\}, q \in \{eq\} \quad (2i)$$

$$\forall p \in \Phi, m, n \quad (2j)$$

$$\forall p \in \Phi, m \quad (2k)$$

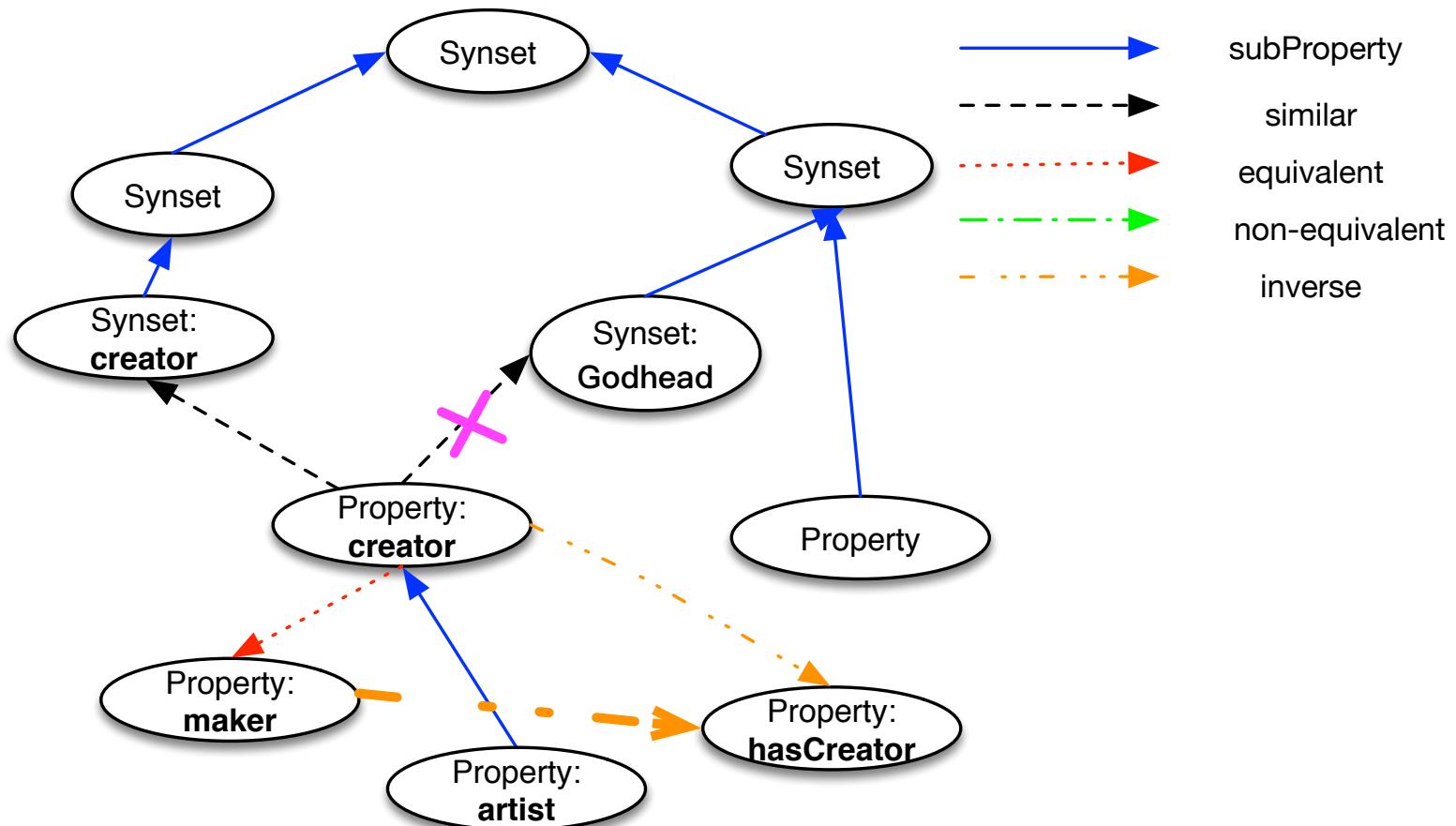
$$\forall p \in \{hasE0P1, hasE0P3, hasE0P5\}, m, n \quad (2l)$$

$$\alpha, C \in \mathbb{R} \quad (2m)$$

$$\Phi \in \{eq, frameEq, nonEq, inverseOf, hasE0P1, hasE0P3, hasE0P5, synFrameEq, subOf\} \quad (2n)$$

$$\Phi \in \{eSyn\} \quad (2o)$$

# Example: Output Graph



# Precision of New Relations

- 8,564 equivalent relations, 1,138 inverse relations, and 2,237 subProperty relations, denoted as Dataset R1, R2, and R3 respectively.

	Precision of R1S	Precision of R2S	Precision of R3S
Annotator 1	0.92	0.77	0.93
Annotator 2	0.89	0.77	0.93

Table 5.9: Precision of output.

- Non-equivalent

Property: <http://purl.org/ontology/daia/availableFor>  
<http://purl.org/ontology/daia/unavailableFor>

Property: <http://purl.org/spar/fabio/isDisciplineOf>  
<http://www.w3.org/2004/02/skos/core#inScheme>

Property: <http://purl.org/spar/fabio/hasDiscipline>  
<http://purl.org/spar/fabio/isSchemeOf>

- Inverse

Property: <http://d-nb.info/standards/elementset/agrelon.owl#hasMuse>  
<http://d-nb.info/standards/elementset/agrelon.owl#isMuseOf>

Property: <http://d-nb.info/standards/elementset/agrelon.owl#hasParent>  
<http://d-nb.info/standards/elementset/agrelon.owl#hasChild>

Property: <http://d-nb.info/standards/elementset/agrelon.owl#hasMurderer>  
<http://d-nb.info/standards/elementset/agrelon.owl#hasMurderree>

Property: <http://d-nb.info/standards/elementset/agrelon.owl#hasStudent>  
<http://d-nb.info/standards/elementset/agrelon.owl#hasTeacher>

Property: <http://d-nb.info/standards/elementset/agrelon.owl#hasEmployer>  
<http://d-nb.info/standards/elementset/agrelon.owl#hasEmployee>

# Our Contribution

- Overcoming **the heterogeneity of vocabularies** on the Web of Data.
- **Enabling new knowledge-driven applications** such as question answering across large numbers of datasets.

# Overall Contributions and Significances

- PEAK: The first method to automatically assess summary content using the pyramid method that also generates the pyramid content models.
- HiText: A novel approach for supporting the reading process.
- LOVE: A novel approach to disambiguate and align properties across vocabularies from the Web of Data.

# References

- **HiText: Text Reading with Dynamic Salience Marking.**  
**Qian Yang**, Yong Cheng, Sen Wang, Gerard de Melo.  
*WWW*, 2017.
- **Neural Machine Translation with Pivot Languages.**  
Yong Cheng, **Qian Yang**, Yang Liu, Maosong Sun, Wei Xu.  
*IJCAI*, 2017.
- **Wise Crowd Content Assessment and Educational Rubrics.**  
Rebecca J. Passonneau, Ananya Poddar, Gaurav Gite, Alisa Krivokapic, **Qian Yang**, Dolores Perin.  
*International Journal of Artificial Intelligence in Education*, 10.1007/s40593-016-0128-6, 2016.
- **PEAK: Pyramid Evaluation via Automated Knowledge Extraction.**  
**Qian Yang**, Rebecca J. Passonneau, Gerard de Melo.  
*AAAI*, 2016.
- **Research on the Key Management for AMI System Based on IEC62056 Standard.**  
Jianfeng Yin, Aibo Shi, **Qian Yang**, Xiaohang Zhang, Liusheng Huang, Haitao Ji.  
*Electrical Measurement and Instrumentation*, Issue No. 2, 2016.
- **Semi-loss-tolerant Strong Quantum Coin-flipping Protocol Using Quantum Non-demolition Measurement.**  
**Qian Yang**, Jiajun Ma, Fenzhuo Guo, Qiaoyan Wen.  
*Quantum Information Processing*, 10.1007/s11128-014-0747-5, 2014.

The data and code are available at  
<http://www.larayang.com>.

Thank you!



# References

- Chin-Yew Lin and Eduard Hovy. ROUGE: A package for automatic evaluation of summaries. In Proceedings of the ACL Text Summarization Branches Out Workshop, pages 74–81, 2004.
- Linking Open Data cloud diagram 2014, by Max Schmachtenberg, Christian Bizer, Anja Jentzsch and Richard Cyganiak. <http://lod-cloud.net/>
- Linking Open Data cloud diagram 2017, by Andrejs Abele, John P. McCrae, Paul Buitelaar, Anja Jentzsch and Richard Cyganiak. <http://lod-cloud.net/>
- Ani Nenkova and Rebecca Passonneau. Evaluating content selection in summarization: The pyramid method. In Daniel Marcu Susan Dumais and Salim Roukos, editors, *HLT-NAACL 2004: Main Proceedings*, pages 145–152, Boston, Massachusetts, USA, May 2 - May 7 2004. Association for Computational Linguistics.

# References

- Luciano Del Corro and Rainer Gemulla. Clausie: clause-based open information extraction. In *Proceedings of the 22nd international conference on World Wide Web*, pages 355–366. ACM, 2013.
- Mohammad Taher Pilehvar, David Jurgens, and Roberto Navigli. Align, disambiguate and walk: A unified approach for measuring semantic similarity. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1341–1351, Sofia, Bulgaria, August 2013. Association for Computational Linguistics.
- Daniel Jacob Gillick. *The Elements of Automatic Summarization*. PhD thesis, EECS Department, University of California, Berkeley, May 2011.
- David N. Rapp and Paul Broek. Dynamic text comprehension: An integrative view of reading. *Current Directions in Psychological Science*, 14(5):276–279, 2005.

# References

- Patricia L. Carrell. Facilitating esl reading by teaching text structure. *TESOL Quarterly*, 19(4):727–752, 1985.
- Ji Soo Yi. QnDReview: Read 100 CHI papers in 7 hours. In *CHI '14 Extended Abstracts*, pages 805–814. ACM, 2014.
- Alan J. Wecker, Joel Lanir, Osnat Mokryn, Einat Minkov, and Tsvi Kuflik. Semantize: Visualizing the sentiment of individual document. In Proc. AVI 2014, pages 385–386. ACM, 2014.
- Valeria de Paiva, Daírio Oliveira, Suemi Higuchi, Alexandre Rademaker, and Gerard de Melo. Exploratory information extraction from a historical dictionary. In Proc. Workshop on Digital Humanities and e-Science at the 10th IEEE International Conference on e-Science, 2014.

# References

- Ed H. Chi, Lichan Hong, Michelle Gumbrecht, and Stuart K. Card. ScentHighlights: Highlighting conceptually-related sentences during reading. In Proc. IUI, 2005.
- Ryan Kiros, Yukun Zhu, Ruslan Salakhutdinov, Richard S. Zemel, Antonio Torra Iba, Raquel Urtasun, and Sanja Fidler. Skip-thought vectors. CoRR, abs/1506.06726, 2015.
- Gü̈nes Erkan and Dragomir R Radev. Lexrank: Graph-based lexical centrality as salience in text summarization. Journal of Artificial Intelligence Research, 22:457–479, 2004.
- Ani Nenkova and Lucy Vanderwende. The impact of frequency on summarization. Microsoft Research, Redmond, Washington, Tech. Rep. MSR-TR-2005-101, 2005.
- H. P. Luhn. The automatic creation of literature abstracts. IBM Journal of Research and Development, 2(2):159–165, 1958.
- George Miller and Christiane Fellbaum. WordNet: An electronic lexical database, 1998.