

# Advancing Natural Language Understanding with Collaboratively Generated Content

#### **Evgeniy Gabrilovich**

Yahoo! Research gabr@yahoo-inc.com

NAE Frontiers of Engineering Symposium, September 2011

Making computers understand human language (at least to some extent ⓒ)

- The importance of world knowledge for natural language understanding
- Paradigm shift: from small-scale resources to collaboratively generated content (CGC)



- Distilling world knowledge from CGC repositories
- Future research directions
  - Obtaining semantic clues from behavioral information (e.g., the authoring process)

## So what differentiates us from computers ? dinner



#### Motivating example



### Sources of knowledge

5	Collaboratively Generated Content					
Before			After	T in the		
Britanr	e di la Nica	articles: 65,000 Since 1768	<b>3,600,000</b> articles Since 2001	Image: State of the state o		
Word	Net	entries: 150,000 Since 1985	2,500,000 entries Since 2002	WIKTIONARY the free dictionary		
CYC		assertions: <b>4,600,000</b> Since 1984	1,000,000,000 Q&A Since 2005	YAHOO!. ANSWERS		
			4,000,000,000 image Since 2004	<sup>s</sup> flickr		

#### **Example problem:**

Semantic relatedness of words and texts

cat ⇔ mouse

How related are



Augean stables ⇔ golden apples of the Hesperides 慕田峪 (Mutianyu) ⇔ 司马台 (Simatai)

#### Used in many applications

- Information retrieval
- Word-sense disambiguation
- Error correction ("Web site" or "Web sight")

#### Semantic relatedness using Wikipedia

#### Structure (links)

- Strube & Ponzetto '06 (WikiRelate)
- Milne & Witten '08 (**WLM**)
- Content (text)
  - $\bigcirc$
- Gabrilovich & Markovitch '09 (ESA)
- □ Structure + content
- Yeh et al. '09 (WikiWalk)

Using Wikipedia structure

### WikiRelate [Strube & Ponzetto '06]

#### □ Find Wikipedia articles titled with the given words

- Attempt joint disambiguation, if possible
  - King  $\rightarrow$  Monarch, King (chess), King (magazine), ...
  - Rook → Rook (bird), <u>Rook (chess)</u>, Rook (card game), …

Otherwise, choose the first sense

#### Compute relatedness

Compare full article texts or use the category hierarchy





Using Wikipedia structure

#### Wikipedia Link-based Measure (continued) Not all links are equal !

#### Loosely related



#### Science

From Wikipedia, the free encyclopedia

This article is about the general term, particularly For other uses, see Science (disambiguation).

Science (from Latin: scientia meaning "knowledge") is that of Aristotle, for whom scientific knowledge was

#### Very general concept

#### Tightly related



#### Atmospheric thermodynamics

From Wikipedia, the free encyclopedia

Atmospheric thermodynamics is the study of heat to work moist air, formation of clouds, atmospheric convection, bounc convection parameterizations in numerical weather models, a





Using Wikipedia content

#### The circular dependency problem







#### Every Wikipedia article represents a concept

14



Leopard

Using Wikipedia content

#### Wikipedia can be viewed as an ontology -

a collection of concepts



Using Wikipedia content

# The semantics of a word is a vector of its associations with Wikipedia concepts



#### Concept and word representation in Explicit Semantic Analysis [Gabrilovich & Markovitch '09]

Panthera	Cat [0.92]	Leopard [0.84]	Ro	bar [0.77]				
Panthera	/							
From Wikipedia, the free encyclopedia /								
Not to be confused with Pantera.								
For other uses, see Panthera (disambiguation).								
Panthera is a genus of the family	Felidae (cats) which contains four	well-known living species: the tiger, the result of the second subfamily. The big cats the second second second	the lion, the	Panthera <sup>[1]</sup>				
while technically referring to all members of the genus, is commonly used to specifically designate the black panther.								
Only the four Panthera cat species	s have the anatomical structure the	at enables them to roar. The primary i	reason for this	Recent				
was formerly assumed to be the incomplete ossification of the hyoid bone. However, new studies show that the ability to								
roar is due to other morphological features, especially of the larynx. The snow leopard, Uncia uncia, which is sometimes								
included within <i>Panthera</i> , does not roar. Although it has an incomplete ossification of the hyoid bone, it lacks the								
special morphology of the larynx. <sup>123</sup>								
Contents [show]								
Name			[edit]	and the state of				

See also: Panther (legendary creature)

According to the *American Heritage Dictionary*, the origin of the word is unknown. A folk etymology derives the word from the Greek πάν *pan*- ("all") and *thēr* P ("beast of prey") because they can hunt and kill almost everything. The Greek word πάνθηρ P, *pánthēr*, referred to all spotted *Felidae* generically. Although it came into English through the classical

#### Computing semantic relatedness as **cosine** between two vectors



#### Experimental results (individual words)



# Augmenting Explicit Semantic Analysis with the **time dimension**

20



#### Temporal Semantic Analysis [Radinsky et al. '11]

21

#### **<u>1. Novel temporal representation of text</u>**



#### 2. Novel temporal text-similarity measurement



Method for computing semantic relatedness using the temporal representation Concept-based information retrieval with Explicit Semantic Analysis [Gabrilovich & Markovitch' 09; Egozi et al. '11]



- Comparing words / short texts → use concepts
- Comparing longer texts use the bag of words & concepts
- Text categorization: assigning category labels to documents
  - Email filtering, news routing, classifying Web pages
- □ Information retrieval: find documents relevant to a query
  - Searching document collections (e.g., Web search)



#### Cross-lingual Explicit Semantic Analysis (CL-ESA)

[Potthast et al., '08; Sorg & Cimiano '08]



## Summary

- Exogenous knowledge is crucial for language understanding
- Collaboratively generated content contains vast amounts of world knowledge
- Concept-based representation addresses two key problems: synonymy and polysemy
- Numerous other resources await to be tapped!
- Future work: Obtaining semantic clues from behavioral information (e.g., the authoring process)

delicious social bookmarking



AHOO! ANSWERS











flickr

# Thank you!

#### gabr@yahoo-inc.com

http://research.yahoo.com/~gabr