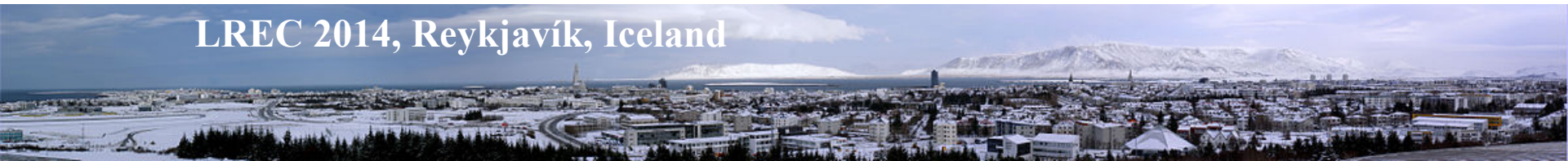




# Thematic Cohesion: Measuring Terms Discriminatory Power Toward Themes

Clément de Groc  
Xavier Tannier  
Claude de Loupy

LREC 2014, Reykjavík, Iceland



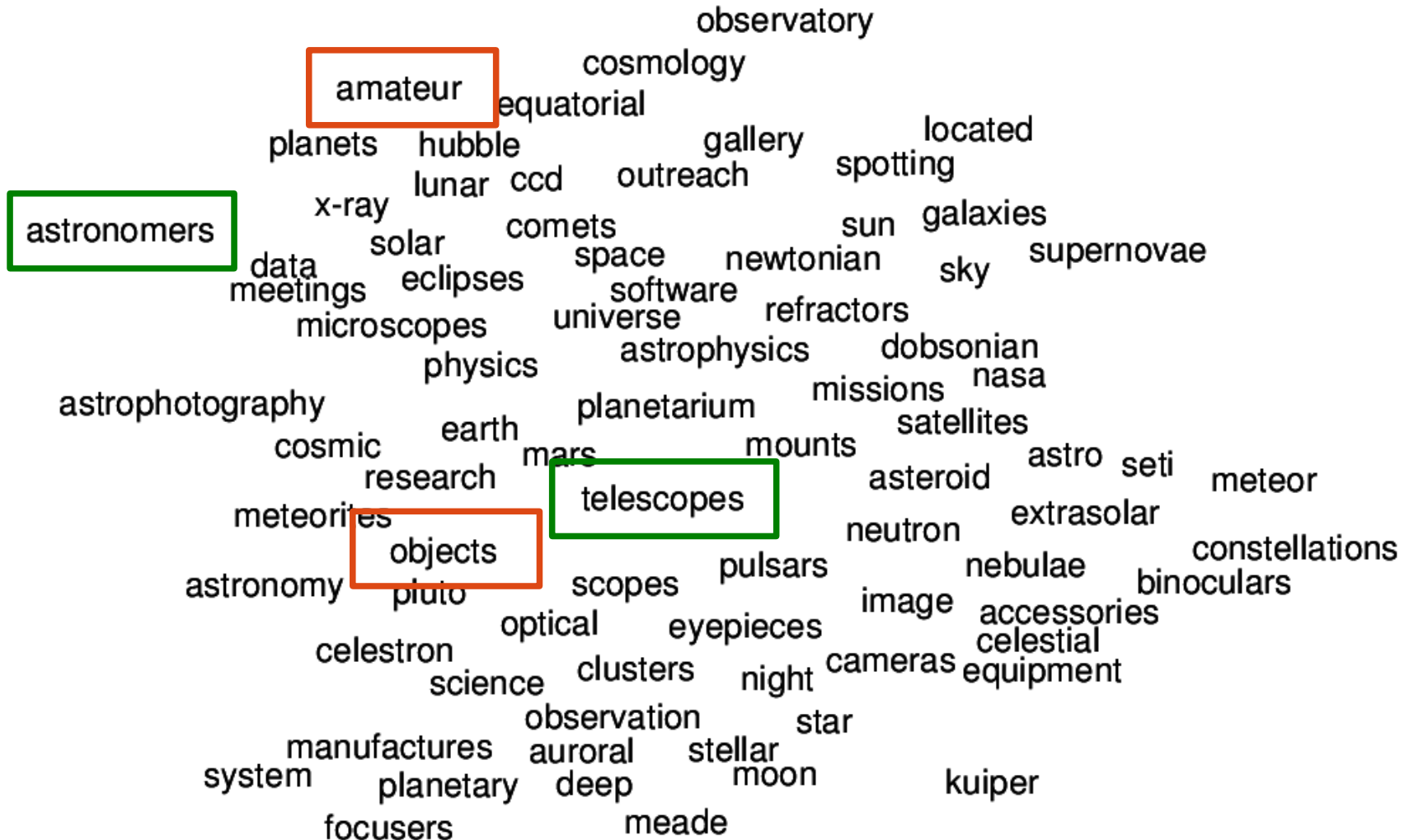
# *Thematic lexicons*

Domain-specific lists of terms (*e.g.* Astronomy)

A word cloud of astronomy-related terms. The words are arranged in a roughly circular pattern, with 'observatory' at the top and 'meade' at the bottom. The terms include: observatory, cosmology, amateur, equatorial, planets, hubble, gallery, located, x-ray, lunar, ccd, outreach, spotting, astronomers, solar, comets, sun, galaxies, data, eclipses, space, newtonian, sky, supernovae, meetings, microscopes, universe, software, refractors, physics, astrophysics, dobsonian, nasa, astrophotography, planetarium, missions, satellites, cosmic, earth, mars, mounts, asteroid, astro, seti, meteor, research, telescopes, neutron, extrasolar, meteorites, objects, scopes, pulsars, nebulae, constellations, astronomy, pluto, optical, eyepieces, image, accessories, binoculars, celestron, science, clusters, night, cameras, equipment, observation, star, manufactures, auroral, stellar, moon, kuiper, system, planetary, deep, moon, kuiper, focusers, meade.

# Thematic lexicons

Domain-specific lists of terms (e.g. Astronomy)



# Thematic lexicons

Domain-specific lists of terms (e.g. Astronomy)

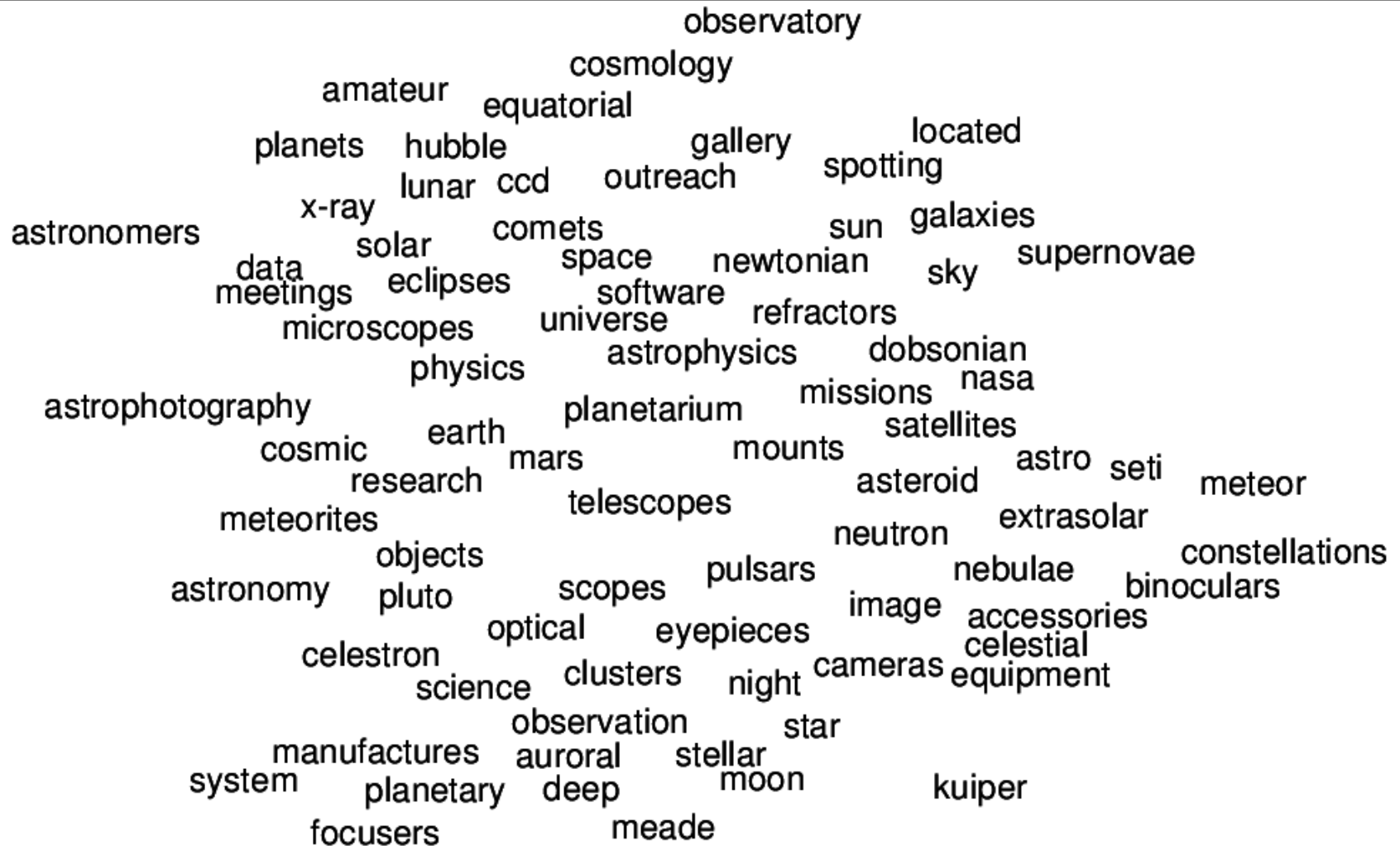
observatory  
amateur equ  
planets hubble  
x-ray lunar cor  
astronomers solar cor  
data eclipses  
meetings microscopes  
astrophotography physics  
cosmic earth m  
research  
meteorites objects  
astronomy pluto opti  
celestron science  
manufactures ok  
system planetary at  
focusers meade

- Building a domain-specific, thematic lexicon is **long** and **expensive**
- Semi-automatic processes exists (Baroni and Bernardini, 2004; Kilgarriff and Grefenstette, 2003; Wang and Cohen, 2007)
- However, a **manual validation** step is still required when high quality resources are needed

# *Objectives*

# Objectives

**Input:** a thematic lexicon



A word cloud of astronomy-related terms, including: observatory, cosmology, amateur, equatorial, planets, hubble, gallery, located, x-ray, lunar, ccd, outreach, spotting, astronomers, solar, comets, sun, galaxies, data, eclipses, space, newtonian, sky, supernovae, meetings, microscopes, universe, software, refractors, physics, astrophysics, dobsonian, nasa, astrophotography, planetarium, missions, satellites, cosmic, earth, mars, mounts, asteroid, astro, seti, meteor, research, telescopes, neutron, extrasolar, constellations, meteorites, objects, scopes, pulsars, image, nebulae, binoculars, astronomy, pluto, optical, eyepieces, cameras, celestial, accessories, celestron, science, clusters, night, equipment, observation, star, manufactures, auroral, stellar, moon, kuiper, system, planetary, deep, meade, and focusers.







# Objectives

**Input:** a thematic lexicon

**Output:** term weights

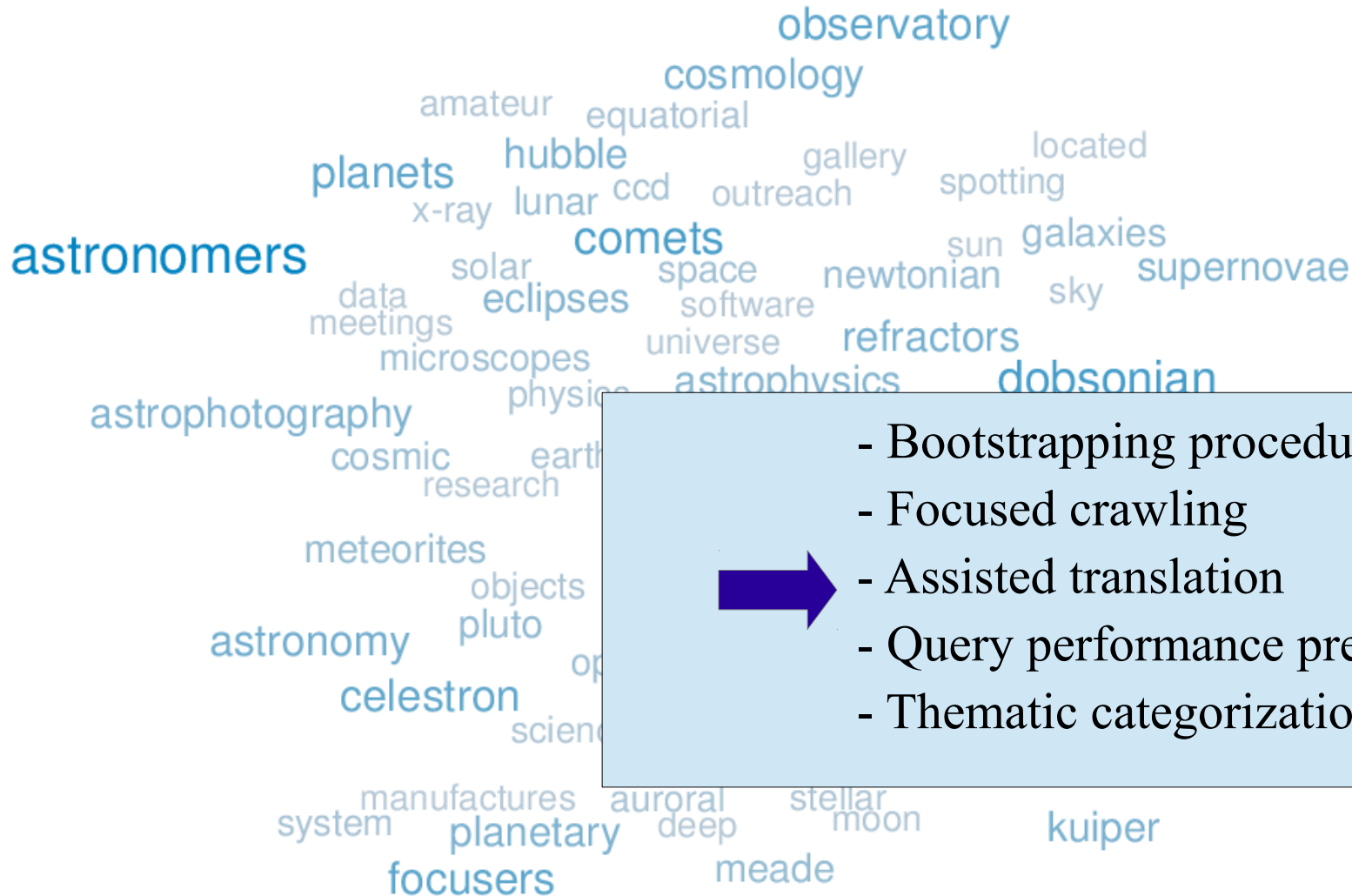
$$L_T = (t_1, t_2, \dots, t_N) \quad \rightarrow \quad w_{L_T} = (w_1, w_2, \dots, w_N)$$

# Objectives



- **Noisy** and **ambiguous** terms are relegated to the bottom of the list (low weights)
- Drastically **reduces manual validation** process
- **Reduces topic drift** through iteration in bootstrapping applications
- Give **better hints** to translators
- ...

# Objectives



# *Thematic Cohesion Value*

# Collecting Data

afterglow  
dwarf stars  
x rays  
red dwarf stars  
accretion disks  
celestial  
coordinates  
bow shocks  
films  
auroral jets  
solar atmosphere  
asteroids  
quasars  
Einstein shift  
space plasmas  
Hubble telescope  
...  
...

Term  
 $t_i$

Standard (general)  
web search engine  
(**blekko**)

$M$  top results

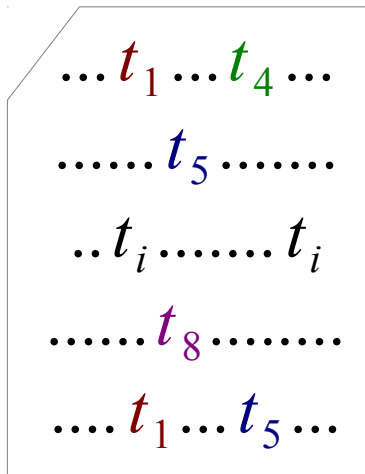
Corpus  $C_i$

... $t_1$ ... $t_4$   
..... $t_5$ ...  
.. $t_i$ ..... $t_i$   
..... $t_8$ .....

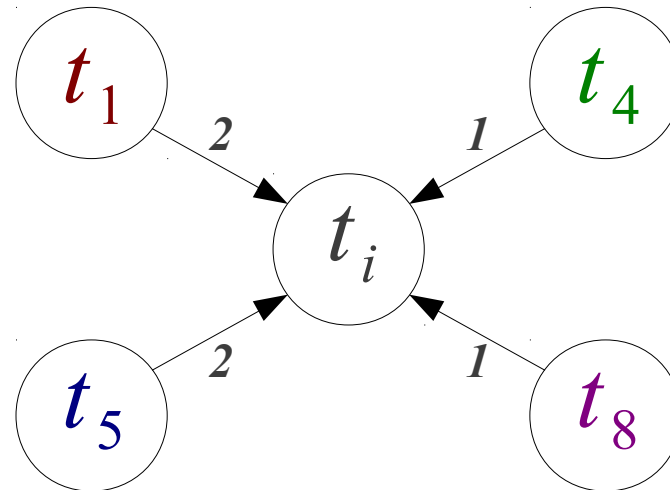
either the **entire web pages**  
or the **search engine**  
**snippets**

# Graph Representation

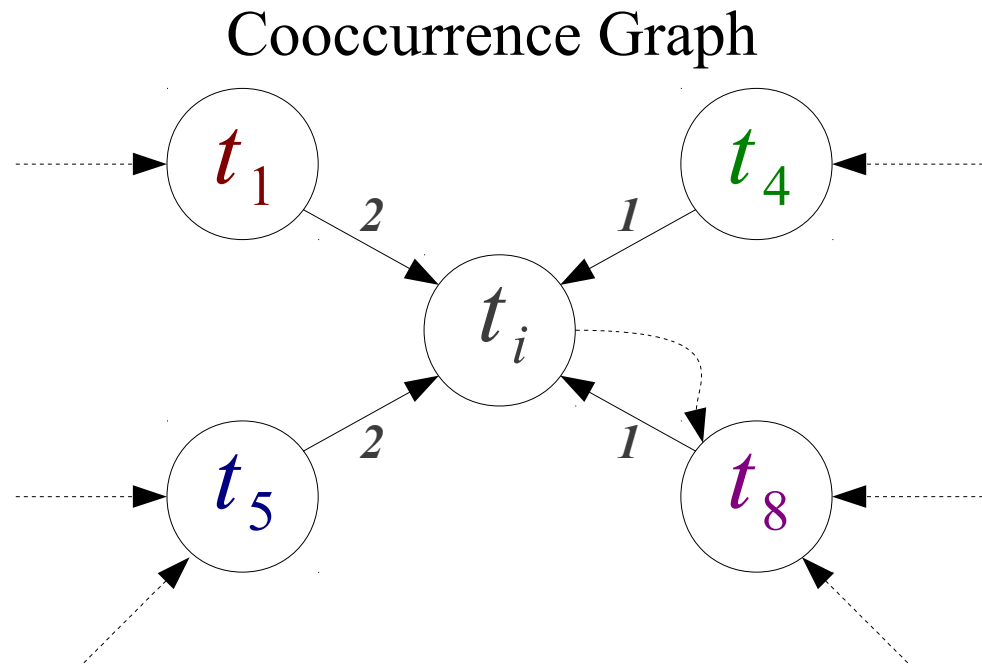
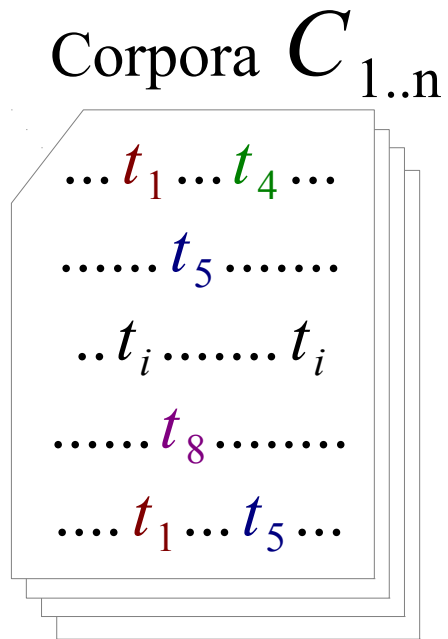
Corpus  $C_i$



Cooccurrence Graph



# Graph Representation

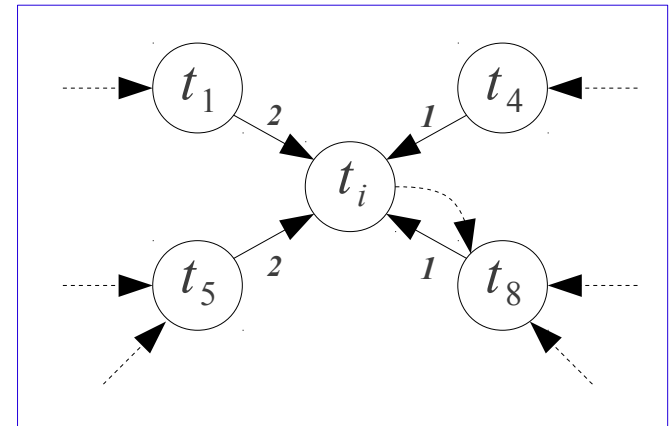


# Thematic Cohesion Value

- **Thematic Cohesion Value** of a word =  
How *central* this word is in the lexicon cooccurrence graph
- Definition #1: In-degree

Score of term  $t_i$  = Number of occurrences of lexicon terms in  $C_i$

$$w_i = \sum_{t_j} n_{t_j, C_i}$$

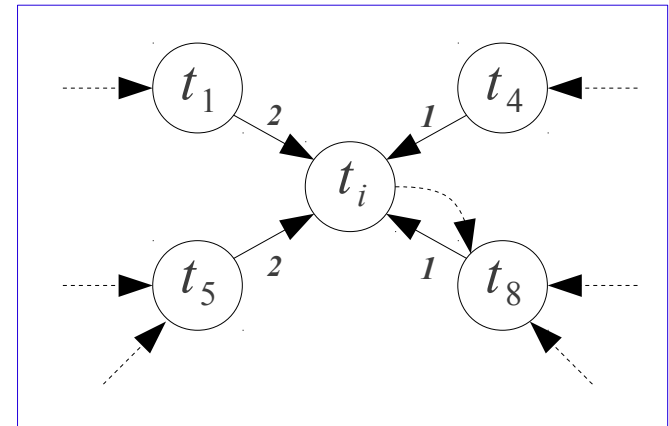




# Thematic Cohesion Value

- **Thematic Cohesion Value** of a word =  
How *central* this word is in the lexicon cooccurrence graph
- Definition #2: Normalized in-degree (score in [0,1])

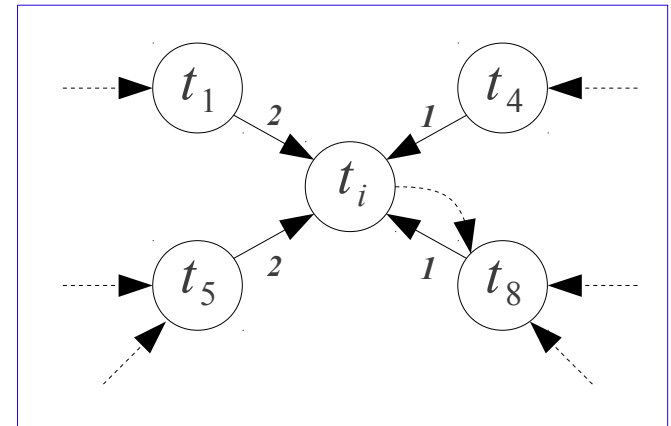
$$w_i = \frac{\sum_{t_j} n_{t_j, C_i}}{\sum_{k \in 1..n} \sum_{t_j} n_{t_j, C_k}}$$



# Thematic Cohesion Value

- **Thematic Cohesion Value** of a word =  
How *central* this word is in the lexicon cooccurrence graph
- Definition #3: Consider other terms weights (**random walk**)

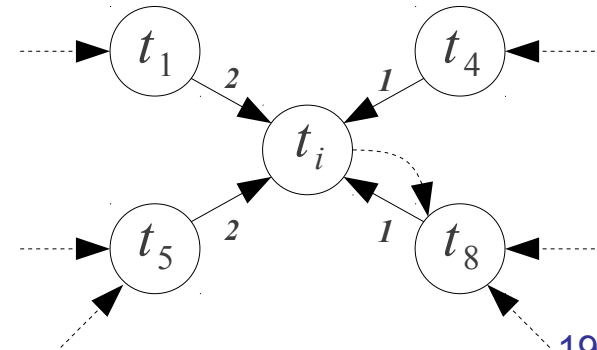
$$w_i = \sum_{t_j} \frac{n_{t_j, C_i} \times w_j}{\sum_{k \in 1..n} n_{t_j, C_k}}$$



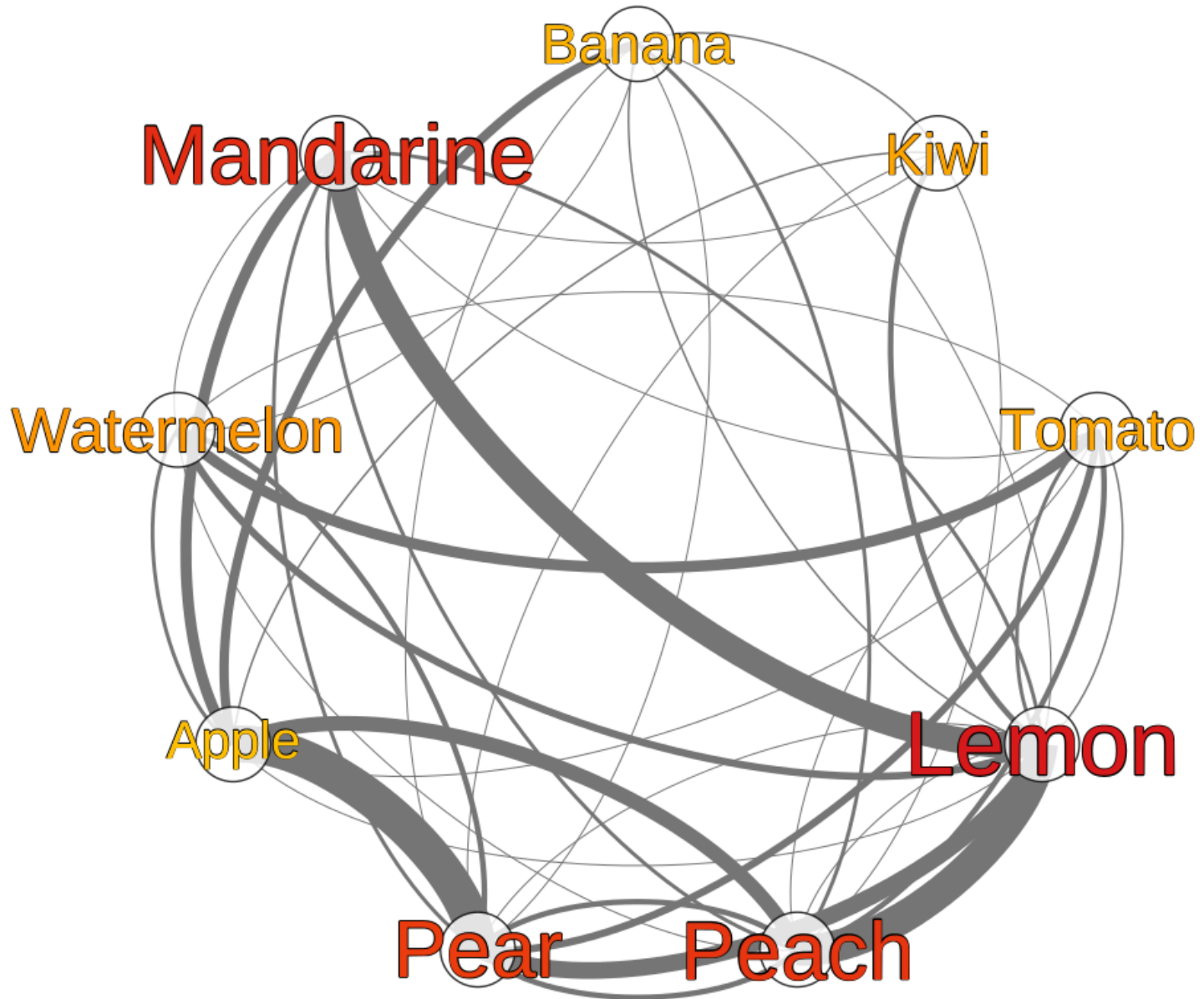
# Thematic Cohesion Value

- **Thematic Cohesion Value** of a word =  
How *central* this word is in the lexicon cooccurrence graph
- Final Definition: Convergence condition → add teleportation vector  
(**PageRank**)

$$w_i = \frac{(1-\alpha)}{N} + \alpha \cdot \sum_{t_j} \frac{n_{t_j, C_i} \times w_j}{\sum_{k \in 1..n} n_{t_j, C_k}}$$



# *Thematic Cohesion Value*



# *Evaluation*

# *Evaluation*

- **Evaluation** of:
  - **Behavior**
    - Influence of the number of documents  $M$
    - Web pages vs. snippets
    - Influence of the initial lexicon size  $N$
  - **Relevance**
    - Best terms should lead to more precise documents for the topic

# *Evaluation : Measure Behavior*

## **Reference lexicons:**

- *Astronomy* (2940 terms, the Astronomy Thesaurus)
- *Statistics* (2752 terms, the ISI Glossary)
- *Medical-1* (2000 terms, MeSH)
- *Medical-2* (2000 terms, MeSH)

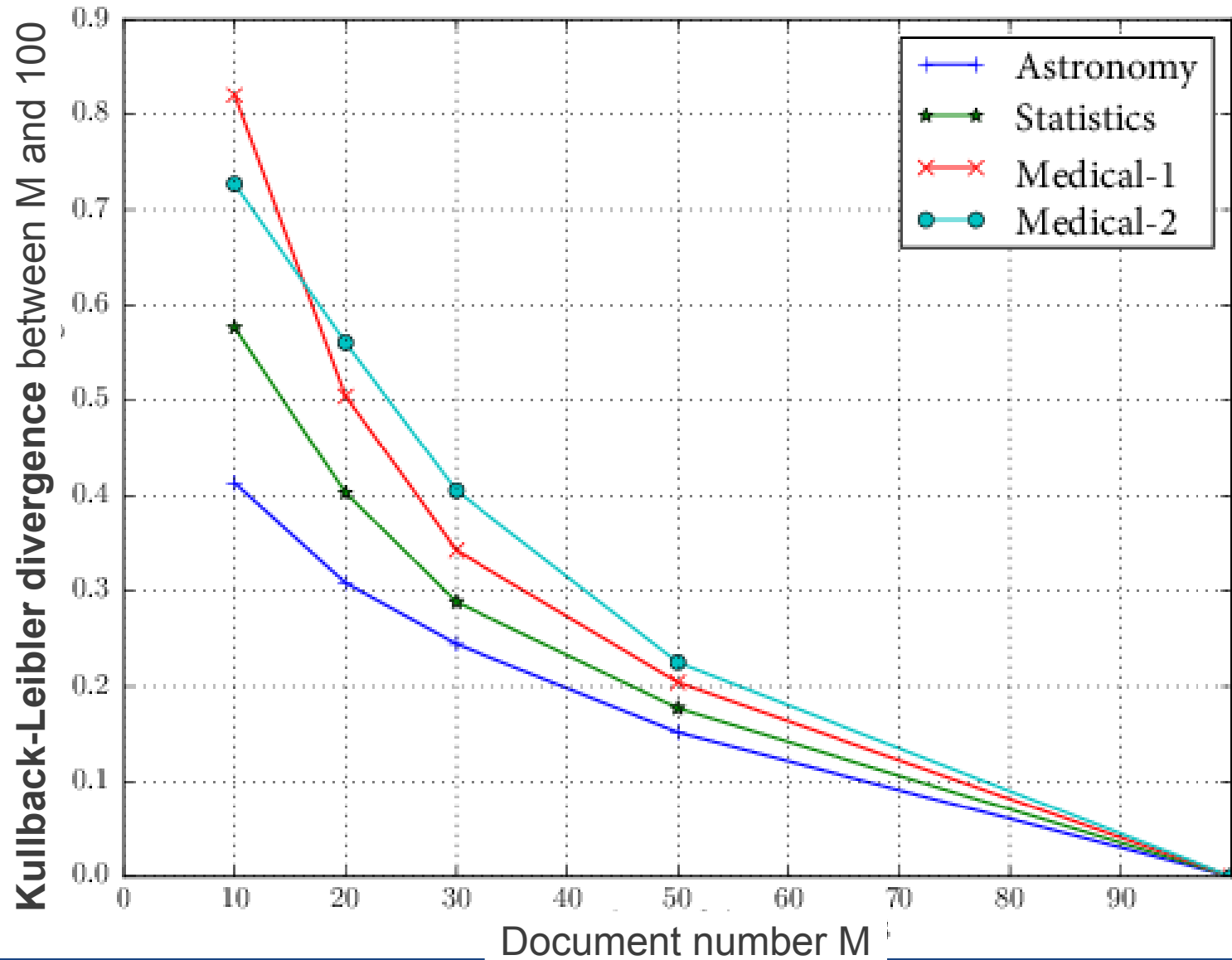
# *Evaluation : Measure Behavior*

- What is the **lowest sufficient number of documents**  
(to obtain stable results) ?
- Study when **results stop changing by adding documents**
  - Different values of  $M$  = number of documents
  - $\max(M) = 100$  web pages or 500 snippets
  - KL-Divergence between lists obtained with  $M$  and  $\max(M)$



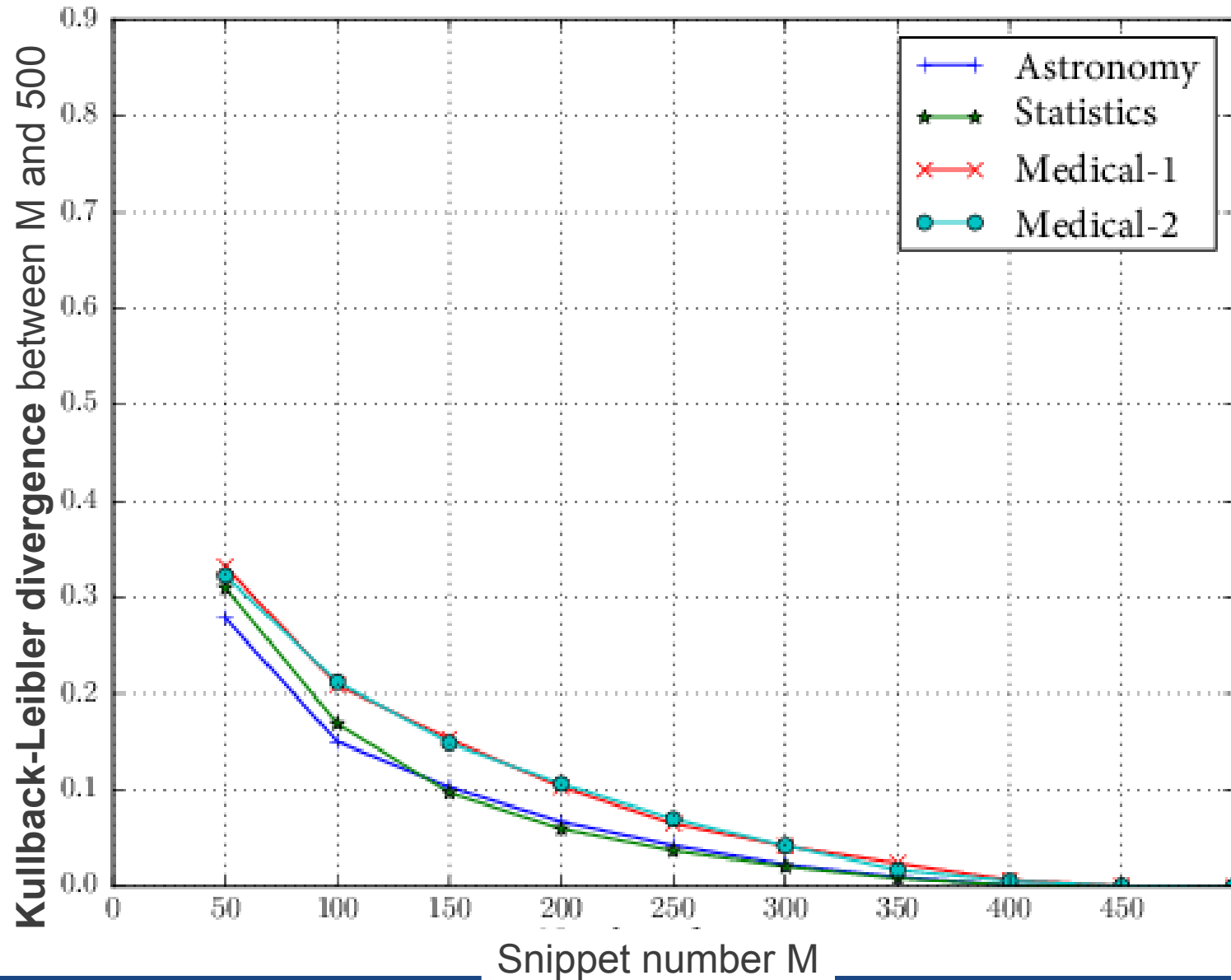
# *Evaluation : Measure Behavior*

Web Pages



# *Evaluation : Measure Behavior*

Snippets

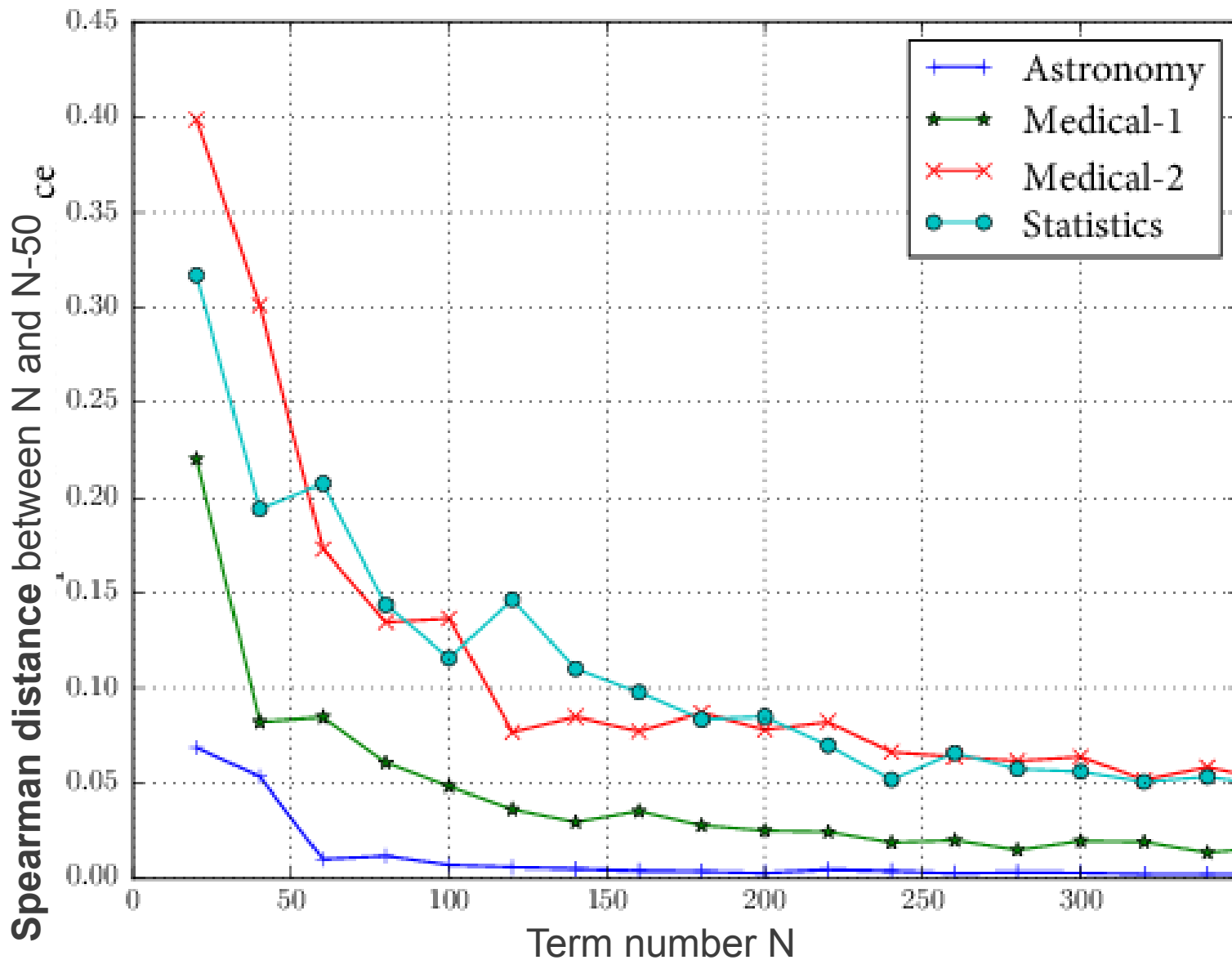


# *Evaluation : Measure Behavior*

- What is the **lowest sufficient number of terms in the lexicon**  
(to obtain stable results) ?
- Study when **ranking stops being messed up by adding new terms**  
to the lexicon
  - Different values of  $N$  = number of terms (between 20 and 1000)
  - Spearman distance of the ranked lists with  $N$  and  $N+50$  terms.

# *Evaluation : Measure Behavior*

Terms



# Evaluation : Relevance

- OpenDirectory (DMOZ)

The screenshot shows the DMOZ website header with the logo, a search bar, and navigation links. Below the search bar, there is a grid of category links. A large arrow labeled 'Categories' points from the bottom left towards the category grid.

**dmoz** In partnership with **Aol.**  
Follow @dmoz | about dmoz | dmoz blog | suggest URL | help | link | editor login

Search [advanced](#)

<a href="#">Arts</a> <a href="#">Movies</a> , <a href="#">Television</a> , <a href="#">Music</a> ...	<a href="#">Business</a> <a href="#">Jobs</a> , <a href="#">Real Estate</a> , <a href="#">Investing</a> ...	<a href="#">Computers</a> <a href="#">Internet</a> , <a href="#">Software</a> , <a href="#">Hardware</a> ...
<a href="#">Games</a> <a href="#">Video Games</a> , <a href="#">RPGs</a> , <a href="#">Gambling</a> ...	<a href="#">Health</a> <a href="#">Fitness</a> , <a href="#">Medicine</a> , <a href="#">Alternative</a> ...	<a href="#">Home</a> <a href="#">Family</a> , <a href="#">Consumers</a> , <a href="#">Cooking</a> ...
<a href="#">Kids and Teens</a> <a href="#">Arts</a> , <a href="#">School Time</a> , <a href="#">Teen Life</a> ...	<a href="#">News</a> <a href="#">Media</a> , <a href="#">Newspapers</a> , <a href="#">Weather</a> ...	<a href="#">Recreation</a> <a href="#">Travel</a> , <a href="#">Food</a> , <a href="#">Outdoors</a> , <a href="#">Humor</a> ...
<a href="#">Reference</a> <a href="#">Maps</a> , <a href="#">Education</a> , <a href="#">Libraries</a> ...	<a href="#">Regional</a> <a href="#">US</a> , <a href="#">Canada</a> , <a href="#">UK</a> , <a href="#">Europe</a> ...	<a href="#">Science</a> <a href="#">Biology</a> , <a href="#">Psychology</a> , <a href="#">Physics</a> ...
<a href="#">Shopping</a> <a href="#">Clothing</a> , <a href="#">Food</a> , <a href="#">Gifts</a> ...	<a href="#">Society</a> <a href="#">People</a> , <a href="#">Religion</a> , <a href="#">Issues</a> ...	<a href="#">Sports</a> <a href="#">Baseball</a> , <a href="#">Soccer</a> , <a href="#">Basketball</a> ...

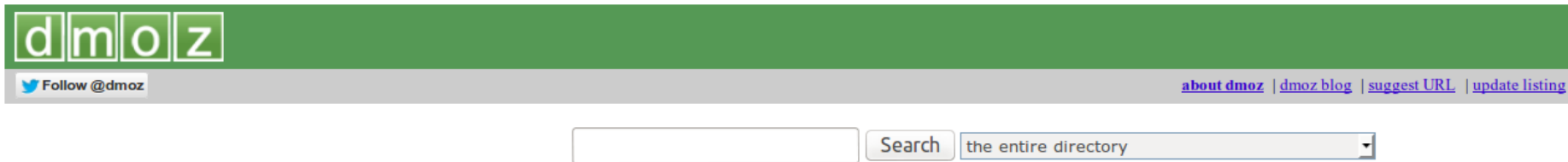
[World](#)  
[Català](#), [Česky](#), [Dansk](#), [Deutsch](#), [Español](#), [Esperanto](#), [Français](#), [Galego](#), [Hrvatski](#), [Italiano](#), [Lietuvių](#), [Magyar](#), [Nederlands](#), [Norsk](#), [Polski](#), [Português](#), [Română](#), [Slovensky](#), [Suomi](#), [Svenska](#), [Türkçe](#), [Български](#), [Ελληνικά](#), [Русский](#), [Українська](#), [العربية](#), [עברית](#), [עברית](#), [תענית](#), [日本語](#), [简体中文](#), [繁體中文](#), ...

[Become an Editor](#) Help build the largest human-edited directory of the web

Copyright © 1998-2014 AOL Inc.

# Evaluation : Relevance

- OpenDirectory (DMOZ)



[Top](#): [Science](#): [Social Sciences](#): [Linguistics](#): [Computational Linguistics](#): [Conferences](#) (33)

- [2000](#) (3)
- [2001](#) (5)
- [2002](#) (3)
- [2003](#) (6)
- [2004](#) (3)
- [2005](#) (0)
- [2006](#) (0)
- [2007](#) (0)
- [2008](#) (0)
- [pre-2000](#) (5)

See also:

- [Computers: Artificial Intelligence: Conferences and Events](#) (52)
- [Computers: Computer Science: Conferences](#) (179)
- [Computers: Human-Computer Interaction: Conferences](#) (38)
- [Science: Social Sciences: Linguistics: Conferences](#) (91)

- [5th international Conference of the Global WordNet Association](#) - The GWC-2010 will be held at IIT Bombay, Mumbai, India, during 31st Jan - 4th Feb, 2010.
- [The Association for Computational Linguistics: Conferences](#) - Information on upcoming ACL (and associated) conferences. Also archives and mirrors of past conferences.
- [Conferences and Workshops](#) - Listed by year, and maintained by the community.

[LREC Conferences](#) - The International Conference on Language Resources and Evaluation is organised by ELRA biennially with the support of institutions and organisations involved in HLT. large number of people working and interested in HLT.

*Site description*



# Evaluation : Relevance

- **Resources:**

- Index sites from **340 categories** of DMOZ (second-level)
- For each category, build a lexicon with **200 best tf.idf terms** from all site descriptions

- **Methodology:**

- Compute **Thematic Cohesion Values** for each lexicon
- Issue each term as a query to our OpenDirectory search engine
- **Idea:** *the more a term is relevant to a topic, the better the precision of retrieved documents*
  - **Compute** the average precision of the set of retrieved documents.
  - **Compare** average precision and thematic cohesion scores  
(Spearman coefficient)

# *Evaluation : Relevance*

Measure	Spearman $\rho$	Significance (% of categories where p-value < 0.05)
<i>tf.idf</i>	0.200	32%
Thematic cohesion	0.434	74%



# *Conclusion*

# *Conclusion*

- A novel Thematic Cohesion Measure
- Weights thematic lexicon terms according to their discriminatory power toward the theme
- Use of a general search engine
- Snippets are more robust and as relevant as web pages
- Useful for corpora bootstrapping, assisted translation, query performance prediction, etc.