# Content Mining,
## a short introduction to practices and policies

**Summary of**

**a study for the Publishing Research Consortium**

**into Journal Article Mining,**

**By Eefke Smit, Maurits van der Graaf (2011)**

**Full study available on PRC website**

# Let's start with a potential user (1)

**Use-case-1:**

**keeping up-to-date**

- Since 1982:
    - 90,000 journal articles on neuroregeneration (e.g. spinal cord injury)

- New articles:
    - on average 22 journal articles per day on neuroregeneration

**Prof. Joost Verhaagen PhD, Netherlands Institute for Neuroscience, Amsterdam**



*stm*

# Let's start with a potential user (2)

## Use-case-2: Information needed as result of laboratory experiments

**Prof. Joost Verhaagen PhD, Netherlands Institute for Neuroscience, Amsterdam**

- Which molecules do play a role in this process?

- Typical outcome of an experiment: hundreds of molecules show enhanced activity

- Next step: how to filter out the relevant molecules?

- *'You would like to have for each of these molecules a meta-analysis about what is already known about these molecules in other processes'*

*stm*

# The essence of TDM is:

**So much information to analyse:**

**Can a machine do this for him ?**

Text mining tool for semantic search by PubTator, see
http://www.ncbi.nlm.nih.gov/CBBresearch/Lu/Demo/PubTator/tutorial/index.html

# Typical text mining consists of

- Processing large corpora of text in an automated way
- To identify entities, instances, actions, relationships and patterns and also for assertion- and sentiment analysis
- For example: genes, proteins, gene-disease patterns, compound properties, chemical structures, side effects of drugs
- **Text mining output typically consists of:**
  - Article clusters and categorisations, indexes
  - Topical maps, to show the occurrence of topics and their inter-relationships
  - Databases with facts, patterns, relationships, statements, assertions, properties found in the articles,
  - Visualisations like graphs, mappings, plot-graphs and topicical maps

# Optimists and Pessimists on TDM

Skeptics:

- Has always over-promised
- Only in specialized fields
- Tools still complicated
- Manual curation necessary
- High investments
- Domain dependent
- No common dictionary
- Overambition in the promise of knowledge discovery

Optimists:

- Vast digital corpus available and growing
- More and more application areas (business, legal, social, etc)
- Tools improving fast
- Manual work reduced
- Public domain or domain precision
- Processing power less of a problem, analytical tools better, visualisation adds to analysis

*stm*

# Study commisioned by
# the Publishing Research Consortium

- Authors:
  - Eefke Smit,
  - Maurits van der Graaf, Pleiade Management & Consultancy
- Two parts:
  - Qualitative study:
    - 29 interviews with experts in academia, research, libraries, vendors and publishers
  - Quantitative study
    - Survey among publishers (members Crossref & STM)
    - 190 responses
- Full report on PRC website www.publishingresearch.net
- Article in the 1st issue of 2012 of Learned Publishing

# Publishers are optimistic:
## Opinions/ expectations for Content Mining in the next 3 years

# Publishers are optimistic, continued:
**Opinions/ expectations for Content Mining on scholarly content in the next 3 years**



Scholarly publishers will mine their content for the purpose of content enrichment, semantic tagging and better navigation.

More new services like Mendeley and Citeseer will emerge as a result of better content mining technology.

The institutional repository world will use content mining for better discoverability of their content.

The investments needed for semantic tagging of scholarly content will be a limiting factor.

The lack of real use cases and proven business benefits is a limiting factor to semantic tagging.

Content mining will remain limited to certain subject fields (such as biomedicine, chemistry) where it was applied first.

■ very much   ■ somewhat   ■ neutral   ■ hardly   ■ not at all

*stm*

# …but publishers do not get many mining requests from 3rd parties:



| | 0 | 10 | 20 | 30 | 40 | 50 | 60 | 70 | 80 | 90 | 100 |
|---|---|---|---|---|---|---|---|---|---|---|---|

Downloading or crawling requests
From corporate customers
From Abstracting and Indexing services
Illegal downloads or crawling
From individual research projects
For derivative information products
For Open Access repositories
Projects with commercial/academic

■ > 10 /year   ■ 5 to 10/y   ■ < 5 /year   ■ never

*stm*

# Publishers are liberal in allowing mining:
## How case-by-case requests are treated



generally to grant permission

to require information about the intent and purpose

to request a financial compensation for commercial purposes

to grant permission for navigational products that drive more traffic to our content

to grant permission for research purposes

to decline for navigational products that compete/replace our content

0  10  20  30  40  50  60  70  80  90  100

■ in 100%   ■ in majority of cases   ■ in some cases   ■ never

*stm*

# …and plan more mining themselves: for retrieval and navigation



| | within the next year | between 1 to 3 years | between 3 to 5 years |

# Cross-sector solutions to facilitate Content Mining better

Suggestions made by experts during the interviews:

1. Standardization of Content Formats
2. One Content Mining platform
3. Commonly agreed access terms
4. One window for mining permissions
5. Collaboration with national libraries

(ad 3: most interviewed experts do NOT see Open Access as a related issue; access terms also relate to datafile delivery or mining on the platform itself)

*stm*

# Survey results for the 5 suggestions for cross-sector solutions
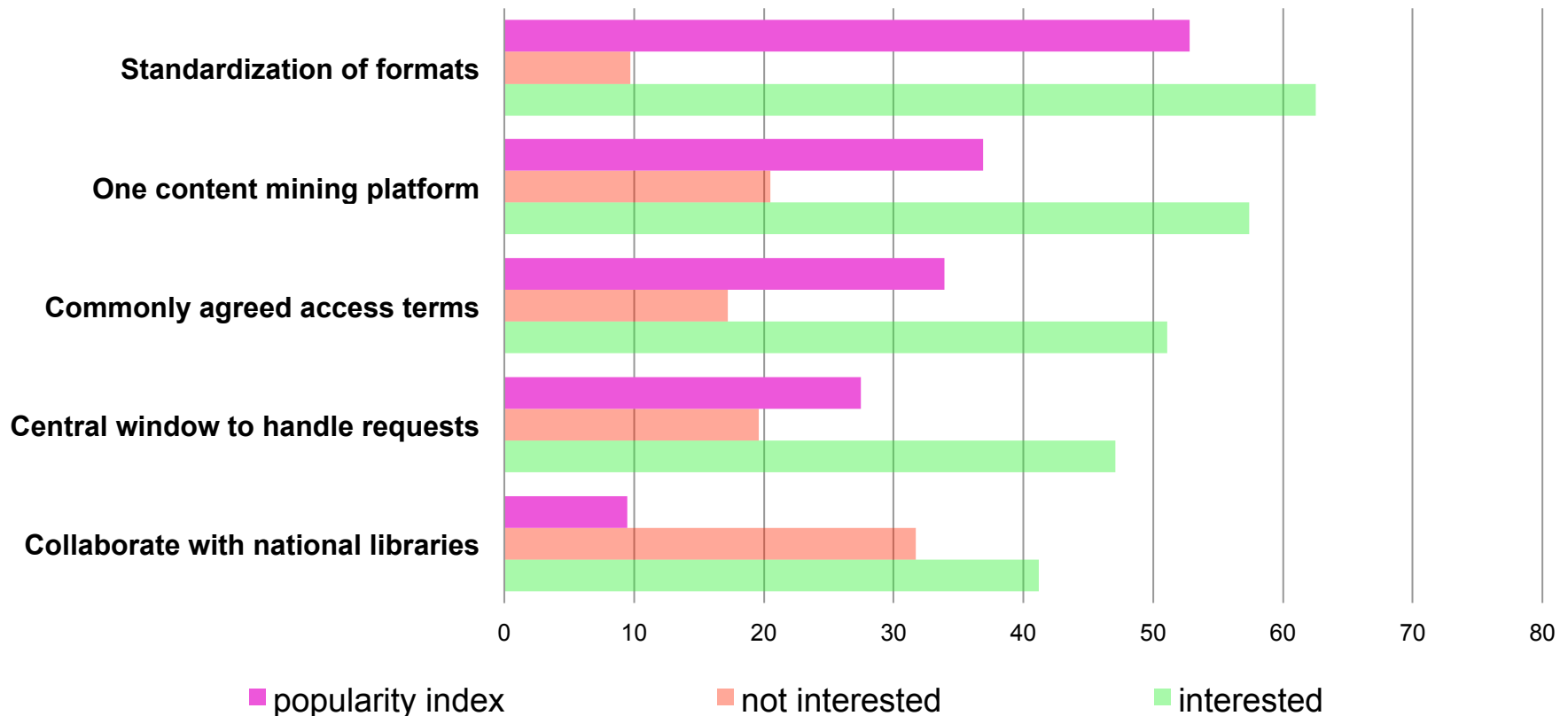## *All respondents*



Chart showing survey results for 5 suggestions across three categories: popularity index (magenta), not interested (salmon), interested (green).

| Suggestion | popularity index | not interested | interested |
|---|---|---|---|
| Standardization of formats | ~53 | ~10 | ~62 |
| One content mining platform | ~37 | ~20 | ~57 |
| Commonly agreed access terms | ~34 | ~17 | ~51 |
| Central window to handle requests | ~27 | ~19 | ~47 |
| Collaborate with national libraries | ~9 | ~32 | ~41 |

■ popularity index   ■ not interested   ■ interested

*stm*

# Survey results for the 5 suggestions for cross-sector solutions: Experts

## *Expert respondents*



**Categories (top to bottom):**
- Standardization of formats
- One content mining platform
- Commonly agreed access terms
- Central window to handle requests
- Collaborate with national libraries

**Legend:** ■ popularity index  ■ not interested  ■ interested

X-axis: 0, 10, 20, 30, 40, 50, 60, 70, 80

*stm*

# Standardisation best prefered, of content formats (and of API's)

Experts believe less in one platform and support standardisation even stronger, not just for content, also for APIs:

Top 3 for all Respondents:
1. Standardisation of Formats
2. One Mining Platform
3. Agreed Permission Terms

Top 3 for Experts only:
1. Standardisation of Formats
2. Agreed Permission Terms
3. One Mining Platform

*stm*

# Questions ?

**Eefke Smit**

**Director Standards and Technology**

**International Association of STM Publishers**

**smit@stm-assoc.org**

*stm*