

# TOP 3

## CHALLENGES FOR COMMERCIAL TEXT MINERS

### How Content Access and Manual Process Hinder Mining Efforts

In biomedical R&D, researchers use text mining tools to extract and interpret facts, assertions and relationships from vast amounts of published information. Mining accelerates the research process, increases discovery and helps companies identify potential safety issues in the drug development pipeline. However, despite the many benefits of text mining, researchers face a number of obstacles before they even get a chance to run queries against the body of biomedical literature. Here are the three primary challenges for researchers as they build a collection of articles (or “corpus”) for their text mining projects:

1

#### INCOMPLETE INFORMATION IN ARTICLE ABSTRACTS

Many researchers build their corpus using scientific article abstracts because they are easily accessible via biomedical databases such as PubMed. While text mining data from abstracts provides some value, there are limitations as to what data can be found within an abstract. The ability to mine the full text of the article – including detailed descriptions of methods and protocols and the complete study results – ensures that researchers don’t miss vital data, discoveries and assertions. However, unlike article abstracts, full text is not often readily available from publishers in a format suitable for text mining.

2

#### LIMITED ACCESS TO XML-FORMATTED CONTENT

When researchers have subscriptions the documents are often provided as PDFs, a format not intended for use with text mining software. Researchers must then spend time converting the PDFs to XML (Extensible Markup Language), the preferred format for use in text mining software. XML is a markup language used to encode documents in a format that is easily read by computers. It is used widely for encoding documents so that computer programs can parse or display the content appropriately. To convert PDFs to XML, researchers must use additional software tools which is not only inefficient but also creates a number of problems with the document itself, including loss of data and tables, conflation of document sections into a “blob of text,” and the addition of bad characters and non-words.

# 3

## INCONSISTENT LICENSING TERMS AND FEES

Because text mining projects depend on access to a broad base of content, businesses must work directly with multiple rightsholders for the use of full-text XML articles, resulting in varying fee structures, inconsistent terms of use and ultimately reduced productivity. Without a common set of terms and conditions for the use of full-text content across publishers, researchers and/or information managers are left with the task of negotiating one-by-one with individual rightsholders to obtain the content and rights they need for text mining.

### LEARN MORE

While text mining can accelerate and enrich your company's R&D program, limited access to full-text content in XML format across multiple publishers and layers of manual process place unnecessary barriers between researchers and the content they want to mine.

Learn how you can improve the results of your text mining efforts, increase efficiency, strengthen your compliance program and save money by contacting RightsDirect at +31-20-312-0437 or emailing [info@rightsdirect.com](mailto:info@rightsdirect.com).

### About RightsDirect

RightsDirect provides licensing solutions that make copyright compliance easy, allowing companies to re-use and share the most relevant digital content across borders. With RightsDirect copyright licenses and complementary information management tools, users can instantly check license coverage, manage permissions and optimize content workflow in one integrated solution.

Based in Amsterdam and with a presence in Tokyo, RightsDirect is a wholly-owned subsidiary of Copyright Clearance Center (CCC). Working in close partnership with the world's leading rightsholders and collecting societies, we offer licensing and content solutions that reflect the needs of local and global organizations. Together, CCC and RightsDirect serve more than 35,000 companies and over 12,000 rightsholders around the globe. [www.rightsdirect.com](http://www.rightsdirect.com)

RD0815

Hoogoorddreef 9  
1101, BA Amsterdam

Tel. +31 20 3120437

[info@rightsdirect.com](mailto:info@rightsdirect.com)  
[www.rightsdirect.com](http://www.rightsdirect.com)

The logo for RightsDirect, featuring the word "RightsDirect" in a blue serif font. Above the letters "i", "g", "h", "t", "s", "D", "i", "r", "e", "c", "t" are several small orange circles of varying sizes, arranged in a slightly curved line.

© 2015 RightsDirect B.V.

A Copyright Clearance Center Subsidiary