# Elsevier's Challenge

Dynamic Knowledge Stores and Machine Translation

Presented By Michelle Gregory, Pascal Coupet

Date 17-05-2016

# OUTLINE

- **Introduction**
  - Elsevier: from publisher to a data & analytics company

- **Elsevier Data**

- **Elsevier Products**

- **Challenges**

- **Current status on Challenges:**
  - Knowledge Graphs
  - Machine aided translation

- **Challenge details:**
  - Creating high quality knowledge graphs
  - Linking taxonomies to translation memory to support machine aided translation

# FROM PUBLISHER TO DATA & ANALYTICS COMPANY

Over the last **50 years** the majority of Noble Laureates have published with Elsevier

Founded over **130 years ago**

Employ over **7,000 employees** in 25 countries

Published over **440,000 articles** in 2015
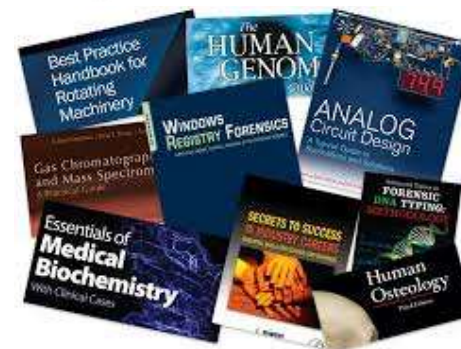
Received over **1.4 million submissions** in 2015

Work with over **30 million** Scientists, students, health & information professionals

Over **61 million** items indexed by Scopus

**CAPABILITIES**

| CONTENT | SOLUTIONS | | | |
|---|---|---|---|---|
| **Elsevier eBooks, Online Journals, Databases** | **Elsevier Research Intelligence** | **Elsevier R+D Solutions** | **Elsevier Clinical Solutions** | **Elsevier Education** |
| Publishes over 2,200 online journals & over 10,000 e-books | Provides universities, governments, and research institutions with the resources and insights to improve institutional research strategy, management, and performance. | Helps corporate researchers, R+D professionals, and engineers improve how they interact with, share, and apply information to solve problems using our digital workflow tools, analytics, and data | Helps medical professionals apply trusted data and sophisticated tools to make better clinical decisions, deliver better care, and produce better healthcare outcomes. | Helps educate highly-skilled, effective healthcare professionals, using the most advanced pedagogical tools and reference works. |
| Cell, ap, THE LANCET, Compendex | Pure, SciVal | Knovel, Geofacets, Embase, Reaxys | ClinicalKey, ToxED | Mosby's Skills+ |

**PLATFORMS**

| ScienceDirect | Scopus | MENDELEY |
|---|---|---|

# ELSEVIER DATA

- **Journals**
  - 3000 journals
  - 440000 articles
  - 1.4 million submissions/year
- **Books**
  - 10000+ eBooks
- **Citations, abstracts and references**
  - 61 million abstracts in Scopus
- **Databases**
  - 26 million substances in Reaxys
  - 4000 drugs in PharmaPendium
  - and more...
- **Taxonomies**
  - 20000 concepts in Omniscience (general subject)
  - 1 million concepts in EMMeT (medicine)
  - 70000 concepts in EmTree (medicine)
  - and more...

# ELSEVIER PRODUCTS

- **Platforms:**
  - ScienceDirect
  - Health Advance
  - Mendeley
- **Products based on analytics:**
  - SciVal
  - Pure
- **Products based on curated data:**
  - Reaxys
  - PharmaPendium
  - Engineering Village
  - Geofacets
  - Pathway Studio

# THE CHALLENGES...

1.  **How to create high-quality non-trivial Knowledge Graphs?**

2.  **Machine Aided Translation:**

    - How to connect/use multi-lingual taxonomies to memory-based translation?
    - How to generate translations of taxonomies?

# STRUCTURED DATA – A COMPETITIVE EDGE

# WHAT WE'VE DONE SO FAR: BUILDING KNOWLEDGE GRAPHS

- **Proof-of-concept work at Elsevier Labs built in 2015**

- **Unsupervised, scalable and built with off-the-shelf technologies**

- **Based on recent work at University College London**

  - Riedel, Sebastian, Limin Yao, Andrew McCallum, and Benjamin M. Marlin. "Relation extraction with matrix factorization and universal schemas." (2013).

# KNOWLEDGE GRAPH - CREATION

- **Elsevier has the data and core structures to fill a knowledge graph –**

  - Semantic Models / Taxonomies

  - Enrichment pipelines
    - Relation and Fact extraction is currently poor but in progress
    - SD Books initiative: extract definitions

  - Open territory:
    - Current glossaries in books
    - Current acknowledgments in books

  - Research Entities: Authors, Institutions, Publications, Journals, …

  - Curated content
    - The right balance between automated processing and hand annotation
    - *Provenance*- proving trusted source can be differentiator for the Elsevier Knowledge Graph

  - Usage Data
    - Co-usage, downloads, popularity ranks

# KNOWLEDGE GRAPH – USES AND APPLICATIONS

- **Flexible disambiguation of entities**

  - Authors, Institutes, Concepts, -- any entity
  - For enrichment pipelines, reference to a knowledge graph with rich data associated with entities will help resolve entities. Enrich entities from:
    - Taxonomies, Wikipedia, DbPedia, Elsevier Sources
  - Powered by existing associations in the graph

- **Query Expansion**

  - Query parsing and interpretation (AskReaxys)
  - Faceting search
    - Suggest associated terms – association of many types (Co-occurrence, taxonomic relations, text-based relations)

- **Recommendations**

  - SD Books use case: background reading
  - Social: often-read together, …

- **Content Generation and Presentation**

  - Question creation
  - Summarization
  - Reasoning: inferred paths (Gene, Physiology, Chemical, Disease)

# CHALLENGE: BEYOND PROOF OF CONCEPT — KNOWLEDGE GRAPHS

- **Construction**
  - What are the productive systems building Knowledge Graph from full-text, full feature articles and patents?
  - What modelling and structuring tooling represents the state-of-the-art in Knowledge Graph creation
  - What evidence is there to show something is state-of-the-art?

- **Valorization**
  - What does the knowledge graph offer that we can't create of higher quality in another way?
  - Ultimate measure is the business value. How can we quantify ROI?
  - What productive instances are there as product offerings – currently – in the space of health, science and technology
  - What could you create to differentiate from the current offerings?

# MACHINE AIDED TRANSLATION

- **Elsevier manually translates all of the assets that need translation:**

  - Books
  - Medical References
  - Clinical Products

- **Problems:**

  - The costs of translation is inhibitive
  - The turn-around time for full text translations is huge: 1-2 years.
  - Machine aided translation only goes to a certain point

- **Elsevier owns translated taxonomies, e.g. English-French-Spanish medical taxonomy EMMeT**

- **Challenge:**

  - How can we connect taxonomies to machine aided translation,
  - How much effort is required to link taxonomies to a translation memory.
  - To control consistency of target language terminology

- **Are there off-the-shelf/ specific/generic methods**

  - Generalizable
  - What are the best machine translation offerings that integrate and conform with Elsevier's multilingual assets

- **Are there off-the-shelf taxonomy translation products**

  - Proven in the market