# Text Mining Techniques

John McNaught
Deputy Director
National Centre for Text Mining
John.McNaught@manchester.ac.uk

# Overview

- Gentle introduction to TM

- Applications/case studies
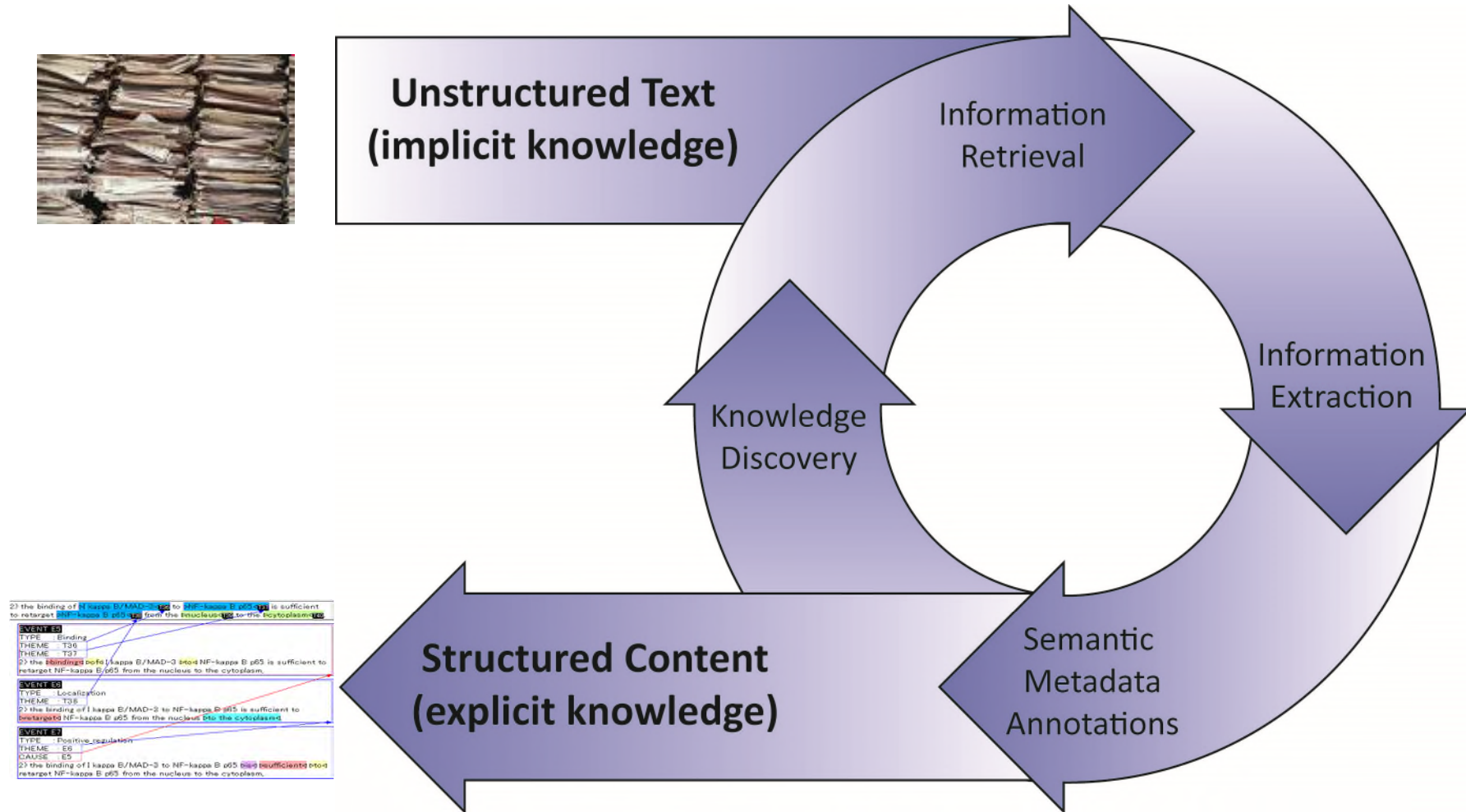  - To help understand what TM could do for you

# NaCTeM

- First publicly-funded (JISC) national text mining centre in the world
  - Provide services to research community
- Initial focus on biology, then social sciences, medicine, chemistry, …
- Processing on a large scale, e.g. for UKPMC (Wellcome T.+17 other funders)
- www.nactem.ac.uk

MANCHESTER
1824
The University of Manchester

NaCTeM
The National Centre for Text Mining

# What is text mining?

- Goal: Discover new knowledge from old
- How:
  - Process (typically) very large amounts of text
    - Millions of documents not unusual
  - Identify and extract information
  - (Link extracted information to already curated knowledge)
  - Mine to discover implicit significant associations
  - Flag (unknown) associations for researcher to investigate further
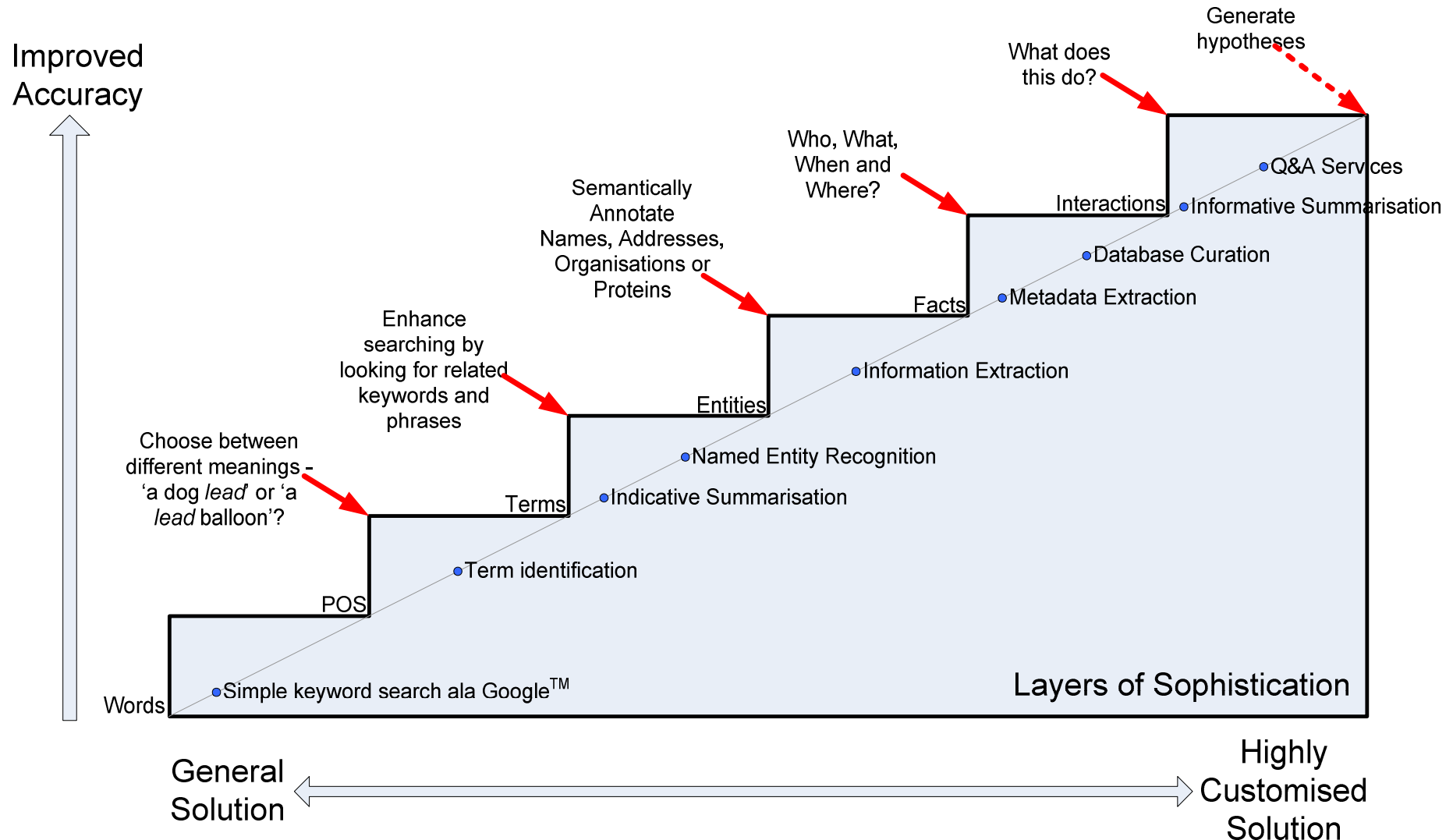  - Spin-off on the way: render information explicit

# From text to new knowledge

# What does TM offer?

- Finds unsuspected knowledge
  - E.g. Disease-gene associations
- Enables discoveries human effort could not achieve (information overload/overlook)
- Enables better search/navigation of literature
  - Semantic search via extracted semantic metadata
- Reduces time spent searching
  - 15-48% of researcher time spent on classic search, 20-50% of classic searches unsatisfied

# A complex problem

- TM involves
  - Many components  (converters, analysers, miners, visualisers, ...)
  - Many resources  (grammars, ontologies, lexicons, terminologies, thesauri, CVs)
  - Many combinations  of components and resources for different applications
  - Many different user requirements  and scenarios, training needs

# Accurate mining of sophisticated knowledge demands customisation

Improved Accuracy

Generate hypotheses

What does this do?

Who, What, When and Where?

Semantically Annotate Names, Addresses, Organisations or Proteins

Enhance searching by looking for related keywords and phrases

Choose between different meanings - 'a dog *lead*' or 'a *lead* balloon'?

- Q&A Services
- Informative Summarisation

Interactions

- Database Curation
- Metadata Extraction

Facts

- Information Extraction

Entities

- Named Entity Recognition
- Indicative Summarisation

Terms

- Term identification

POS

- Simple keyword search ala Google™

Words

Layers of Sophistication

General Solution

Highly Customised Solution

# Customisation

- ## Rule writing
  - – Expert expresses knowledge in formalism
    - Issues of exhaustivity, specificity
- ## Machine learning (supervised)
  - – Expert annotates large amounts of text for desired phenomena
    - Accelerated annotation editor (e.g. NaCTeM's Argo)
      - – System 'completes' annotation as it learns from your decisions
  - – Train machine learning algorithm
    - Apply trained model to find similar patterns of phenomena in new texts
- ## Hybrid approaches
- ## Knowledge poor vs. knowledge rich approaches

# Brief aside

- Many of the systems and applications that follow are based on biomedicine
  - A small intellectual leap required…
- Techniques are portable to other domains and text types
  - Modulo customisation…
  - Modulo understanding user requirements
    - What's a frame (for you, for your colleague)?

# Producing multilevel annotated data: NaCTeM's Argo editor



Extract from McLaughlin et al. (2007) in PLoS Pathogens

# Flexible annotation



Extract from Turk et al. (2007) in Marine Drugs

# Information extraction: many techniques at many levels

# Semantic search using named entity recognition

- Identify named entities in texts
  - People, organisations, locations, genes, drugs, …
- Use entity types in searching
  - Reduce hits to precise, relevant ones
- Make life easy by offering facetted search using the NEs
- NaCTeM's Kleio
- www.nactem.ac.uk/software/Kleio/

# General search

# Add a mined NE facet

# Noisy data

- Full-scale analysis techniques often not suitable for noisy data
  - OCR, digitised texts
  - Emails, unprofessional blogs, tweets, …
- May have to use partial analysis techniques, often based on machine learning
- Even apparently 'clean' data can be hard to process
  - "Hamburger PDFs" (P. Murray-Rust)

# Full-scale analysis

- For formally-written text, full-scale analysis possible
- Deep parsing with MEDIE
  - Allows semantic search
  - Query consists of partial or fully specified fact/event template
  - www.nactem.ac.uk/medie/
  - Developed by Univ. Tokyo (close collaboration)

# MEDIE: semantic search of pre-analysed collection

Sentences from MEDLINE where news is the object, with obesity specified in (here unseen) advanced options

# What users want…

- Many users just give a one-word query when searching
  - Very few explore 'advanced search'
- So NaCTeM uses deep parsing results to *generate questions* relevant to a query that are *known to have answers*
  - UKPMC EvidenceFinder
    - http://labs.ukpmc.ac.uk
  - UKPMC: archive of full texts in life sciences
    - EBI, British Library, Univ of Manchester (Mimas and NaCTeM)

# UKPMC EvidenceFinder

Semantic search based on analysis of full text articles

But currently, EF oriented more towards Biology… and can only text mine OA subset

# UKPMC EvidenceFinder

But asking about e.g. a protein gives better idea of functionality…

330 thousand OA full text articles analysed.

Generated questions known to have answers.

# Finding associations

- Known associations (but not to you)
  - Unknown knowns (the one DR missed out)
- Unknown associations (to anyone)
  - Statistics of surprise
  - A related to B, B to C, so infer A related to C
- www.nactem.ac.uk/services.php
- Check out FACTA+ and FACTA+ visualizer

# Marrying IR with TM

- Classifying
- Clustering
- Summarisation
- Term extraction
- Related documents based on analysis of returned results
- More detail: talk by James Thomas
- Recommendation systems

# Opinion, sentiment mining

- Opinion polling
  - Ratios for/against a proposition
- What ± passions are aroused about
  - A topic?
  - A social group?
- What are the terms in which the debate is conducted?
- How are media used in the formation and dissemination of opinion?

# Case study: Alternative vote debate

- Named entity recognition
  - Finding mentions of political parties, politicians
  - Including alternative designations
- Term and topic discovery
  - Most significant terms used in a (sub) collection of text
- Subjectivity analysis
  - Positive or negative orientation, whether subjective or objective style
- Speech act and rhetorical analysis
  - Whether initiating, reacting, supporting, arguing, citing evidence, etc.
- Semantic search engines based on all the above

# NaCTeM analysis of tweets on AV: overall (document-level) sentiment

# Sentiment analysis of 1 document

# Sentiment: challenges

- Sentiment, subjectivity, emotion
  - Polarity of subjective expression can be confused with being *pro* or *anti* a given proposition ("I support AV" is not subjective)
- Irony
- Rhetoric (last element often overturns predecessors)
- Quality and cultural specificity of resources (e.g. sentiment lexicon)

# Opinion/sentiment mining tutorial

http://www.cs.uic.edu/~liub/FBS/Sentiment-
  Analysis-tutorial-AAAI-2011.pdf

- Recent comprehensive tutorial by Bing Liu
  at AAAI 2011,  August.

# Conclusion

- TM has many aspects, components, levels
  - Different combinations for different tasks
  - Can be applied in any domain (customisation)
- No "one size fits all"
- Domain experts and developers must engage to ensure appropriate application for needs
- Copyright & licensing: support adoption of Hargreaves recommendations!

# Acknowledgements

- Text mining team at NaCTeM (~19)

-  THE UNIVERSITY OF TOKYO

- Funders and sponsors

Wellcome Trust + 17 other sponsors (UKPMC)

Elsevier