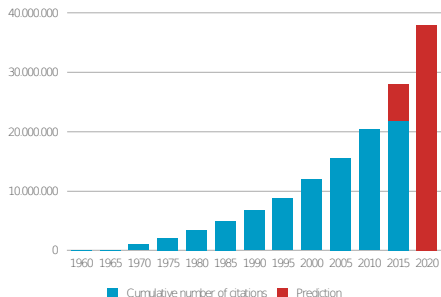# Algorithms for resource-constrained domain-specific knowledge management

### Bachelor's thesis

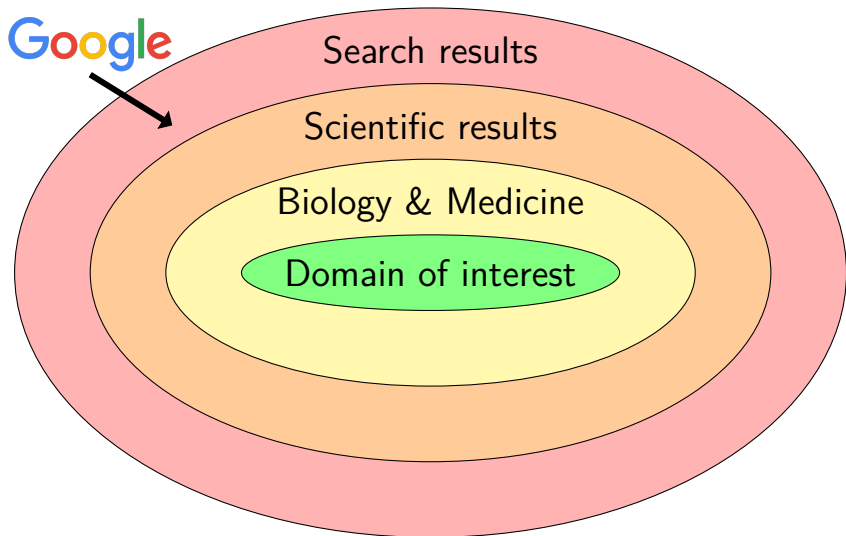Uli Köhler

September 24, 2015

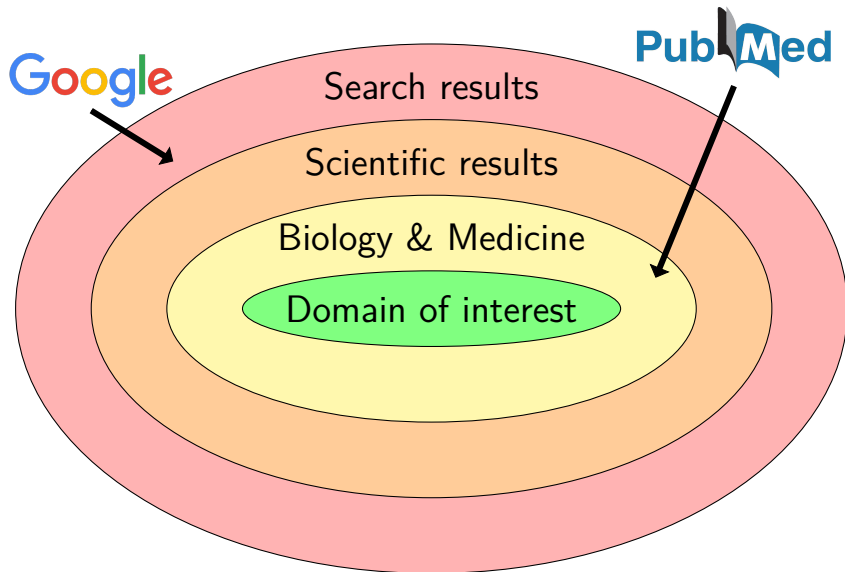# Text mining



- Growing number of publications
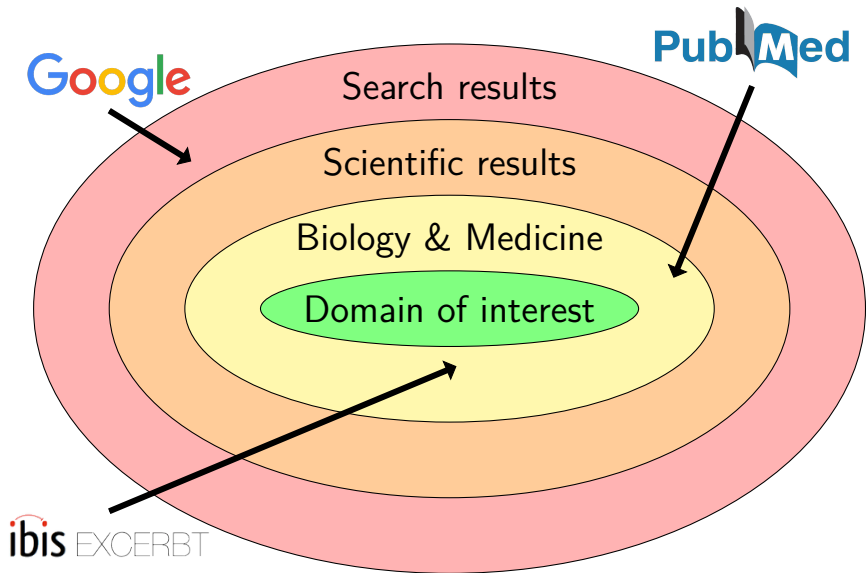
- Quick access to information required

  $\rightarrow$ Text mining

# Search engines

# Search engines

# Search engines

# Domain-specific text mining

- Results from outside of area of interest
  $\rightarrow$ High false-positive rate
- Large-scale text mining is resource-intensive
  $\rightarrow$ Expensive hardware required

- Domain-specific text mining covers only a small area of interest
- Smaller dataset $\rightarrow$ cheap commodity hardware

# TRANSLATRON

- **Transla**tional Bioinformatics **T**ool with **r**elatime **on**tology
- A simple tool for domain-specific text mining
- Easy to use — low hardware requirements
- Web-based user interface
- Real-time search in corpus and ontology
- Named entity recognition (**NER**)

  Example: *The* $\boxed{PrP^{Sc}}$ $\boxed{prion}$ *causes ovine* $\boxed{prion\ diseases}$

# Algorithms for domain-specific text mining

- Conventional algorithms built for large-scale datasets:
  - Hundreds of gigabytes of RAM available
  - Hundreds of terabytes of disk space available
  - Clustered architecture

# Algorithms for domain-specific text mining

- Conventional algorithms built for large-scale datasets:
  - Hundreds of gigabytes of RAM available
  - Hundreds of terabytes of disk space available
  - Clustered architecture

- Novel algorithms required for domain-specific approaches

# Algorithms in *TRANSLATRON*

- *YakDB* High-performance database
- *PRIMORDIAL* text indexing
- *PRAISER* distributed indexing
- *PERSIST* single-token indexing
- *PRESIDE* real-time prefix search
- *PRO-PANE* priority-based result ordering
- ***FiT-NESS* named entity recognition**
- *WESTSIDE* client interface

# FiT-NESS

- **Fi**rst-**T**oken-based **N**amed **E**ntity **S**election **S**cheme
- Trivial: *Single-token entities* like *BRCA1*
- Hard: *Multi-token entities* like *prion diseases*

# FiT-NESS

- **Fi**rst-**T**oken-based **N**amed **E**ntity **S**election **S**cheme
- Trivial: *Single-token entities* like *BRCA1*
- Hard: *Multi-token entities* like *prion diseases*
- *FiT-NESS aproach*:
  - Ignore everything but the first token
  - When we find a hit, check if subsequent tokens match the entity

# *FiT-NESS* II

prion diseases | MeSH:D017096

prion ⟶ prion diseases | MeSH:D017096

# FiT-NESS III

# Key advantages of *TRANSLATRON*

- Can be installed on resource-constrained devices:
  - Notebooks
  - Mobile devices (smartphone, tablet, ...)
  - Embedded devices
- Simple architecture
- Easily adaptable to specific requirements
- Can import internal documents (lab reports, ...)
- Individual installations for each researcher or workgroup

Live demonstration

# Outlook & conclusion

# Outlook & conclusion

- *TRANSLATRON* is only a proof-of-concept
- ... but easier to adapt than conventional tools
- Only basic features are implemented
- Not infinitely scalable
- Applications in disaster relief?
- Applications for text mining with internal documents?

# Acknowledgements

Mathias C. Walter

Prof. Dr. Hans-Werner Mewes

... and many others ...

# Thank you for your attention!

References and sources available at

https://github.com/ulikoehler/Bachelor
https://github.com/ulikoehler/Translatron
https://github.com/ulikoehler/YakDB

Thesis & talk available at http://techoverflow.net
Contact: ukoehler@techoverflow.net

## Questions?

# Image sources

http://www.case.edu/med/nutrition/images/pubmed-logo.jpg

http://mips.helmholtz-muenchen.de/excerbt

http://www.raspberrypi.org/blog/raspberry-pi-2-on-sale/

https://www.raspberrypi.org/blog/raspberry-pi-2-on-sale/

http://www.depts.ttu.edu/hpcc/

Wachinger: Next Generation Knowledge Extraction from Biomedical
Literature with Semantic Big Data Approaches

# Excerbt architecture

# YakDB architecture

# Translatron demo