CompTox

LITERATURE MINING AND KNOWLEDGE DISCOVERY TOOLS FOR VIRTUAL TISSUES Amar V Singh*, Thomas B Knudsen and Imran Shah

U.S EPA, ORD, Computational Toxicology Research Program and *Lockheed Martin, Research Triangle Park, NC.

UNITED STATES ENVIRONMENTAL PROTECTION AGENCY

Science Question

Deriving novel relationships from information in the scientific literature is an important adjunct to data-mining of complex datasets such as ToxCast™ and ToxRefDB. Concept-mining is a knowledge discovery tool to help address causal links, associations, relationships, and patterns among collections of relevant terms and keywords.

The question addressed here pertains to the applications of text-mining algorithms to extract relevant information from the literature and find key relationships.



VT-KB workflow: Implementation to store information such as data, rules and classifications; organize this information into structured concept models (ontologies); and facilitate retrieval of relevant information based on hypothesis-driven queries. The queries may be human readable (manual) or machine-readable (automated).

Methods/Approach





a. co-occurrences of Chemical and DevToxTerms derived from ToxRefDB

b. co-occurrences of DevToxTerms and Genes (GDX) derived from PubMed

c. co-occurrences of Chemicals and Features derived from Text Mining

COMPUTATIONAL

FOXICOLOGY

d. inference for chemical-specific hypothesis when information from a, b and c are combined

Research Goals

- 1. Build text- and data-mining tools to extract and organize relevant facts from the scientific literature and other sources of information.
- 2. Build knowledgebase for Virtual Tissues (VT), in silico models to simulate the cellular fabric of tissues to predict multicellular behaviors in specific biological systems such as the adult liver (v-Liver™) or developing embryo (v-Embryo™)

1-04



Named Entity Recognition (NER) Based Search

Snapshot of @Note Application

Knowledge Acquisition

Hypothesis Generation



research&development



Results/Conclusions

- Text- and data-mining together form a powerful tool for knowledgebase development
- ToxCast 320 chemicals- raw search returned 186K PubMed abstracts covering 82% of chemicals; filtering by <embryo or fetal> - 9K abstracts covering 48% of the chemicals
- Literature derived information when combined with structured data yield novel hypotheses

Impact and Outcomes

- Semantic integration of unstructured, distributed. heterogeneous data to structured data
- Knowledgebase to support generalized, dynamic model development
- Flexible-extensible strategy to support decision-making in risk assessment

Future Directions

- Expand knowledgebase to other domains (eg, vasculogenesis, stem cell biology, alternative models)
- Web-interface to guery the database for literature based discoveries
- Connect to ACToR database for open access

References

- Singh AV, Yang C, Kavlock RJ and Richard AM (2010) Development Toxicology Research Strategies: Toxicology. Computational Comprehensive Toxicology, 2nd edition (Editors Daston GP and Knudsen TB) Elsevier:NY
- Lourenco A et al.@Note: A workbench for Biomedical Text Mining (2009) Journal of Biomedica Informatics Volume 42, Issue 4, Pages 710-720
- Singh Av. Shah I and knudsen TB. (2009) Literature mining and knowledgebase development for v-Embryo (in Preparation)

This poster does not necessarily reflect EPA policy. Mention of trade names or commercial products does not constitute endorsement or recommendation for use