Recent Improvements of the BEL Information Extraction workFlow (BELIEF) for Biomedical Text Mining and Curation

Justyna Szostak¹, Sumit Madan², William Hayes³, Jens Doerpinghaus², Juliane Fluck², Marja Talikka¹, Manuel C. Peitsch¹, Julia Hoeng¹

¹Philip Morris International R&D, Philip Morris Products S.A., Quai Jeanrenaud 5, 2000 Neuchâtel, Switzerland (Part of Philip Morris International group of companies).

²Fraunhofer Institute for Algorithms and Scientific Computing, Schloss Birlinghoven, 53754 Sankt Augustin, Germany

³Applied Dynamics Solutions LLC, Rahway NJ USA

Introduction

Construction of structured knowledge requires technology that links text mining and curation to knowledge repository. We recently presented BEL Information Extraction workFlow (BELIEF) as a tool that facilitates the transformation of unstructured information described in the literature into structured knowledge networks (1). BELIEF automatically captures causal molecular relationships from scientific text and encodes them in BEL statements. BEL (Biological Expression Language) is a computable and human readable language for representing, integrating, storing, and exchanging biological knowledge in causal and non-causal triples. Recently, we have improved the curation process by extending the biomedical vocabulary and by making the curation dashboard more flexible. Moreover, BELIEF was enhanced with the integration of the OpenBEL API that allows direct linkage to the OpenBEL platform and enables upload of curated documents into the BEL knowledge base. These technological developments of BELIEF greatly improve the curation process and make the BEL knowledge more manageable. We continually use the BELIEF to develop an extensively annotated knowledge base of BEL triples that serve as building blocks for causal biological network models.

The BEL Information Extraction WorkFlow (BELIEF) Usecase

ttp://www.openbel.or

NEUROSCIENCE

NEUROSCIENCE

Mentitud punt out it follow directed if the

Science

Maying and
Spinal

Maying and
Spinal

Maying and
M

"Text Mining Pipeline"

Scientific Literature

Cellular and molecular interactions taking place during disease development or progression are written into computable BEL statements

a(CHEBI: oxLDL) -> p(HGNC:TNF)

Exposure

Cellular and molecular interactions

Leucocytes Inflication

Bel Statement:

SUBJECT OBJECT SUBJECT SUBJECT SUBJECT SUBJECT

Tissue level Complication

Atherosclerosis complication

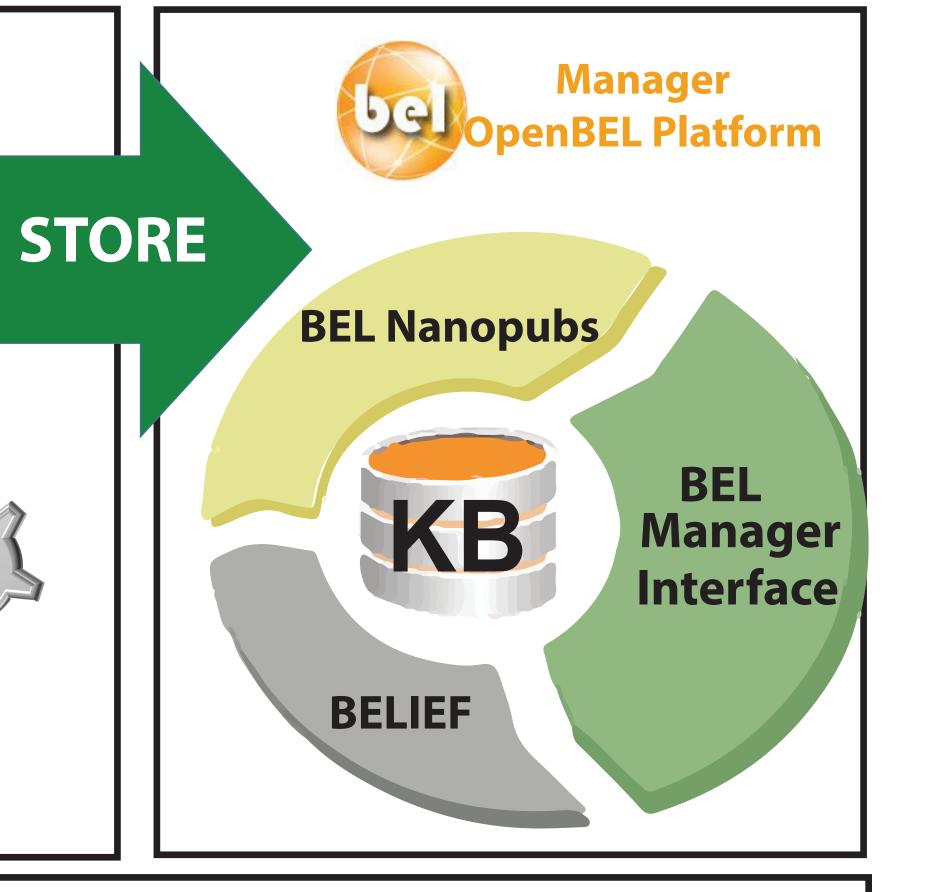
Complication

SUBJECT SUBJECT SUBJECT OBJECT

 $act(p(HGNC:TNF)) \longrightarrow act(p(HGNC:MMP))$

Structured Knowledge

Knowledge Repository

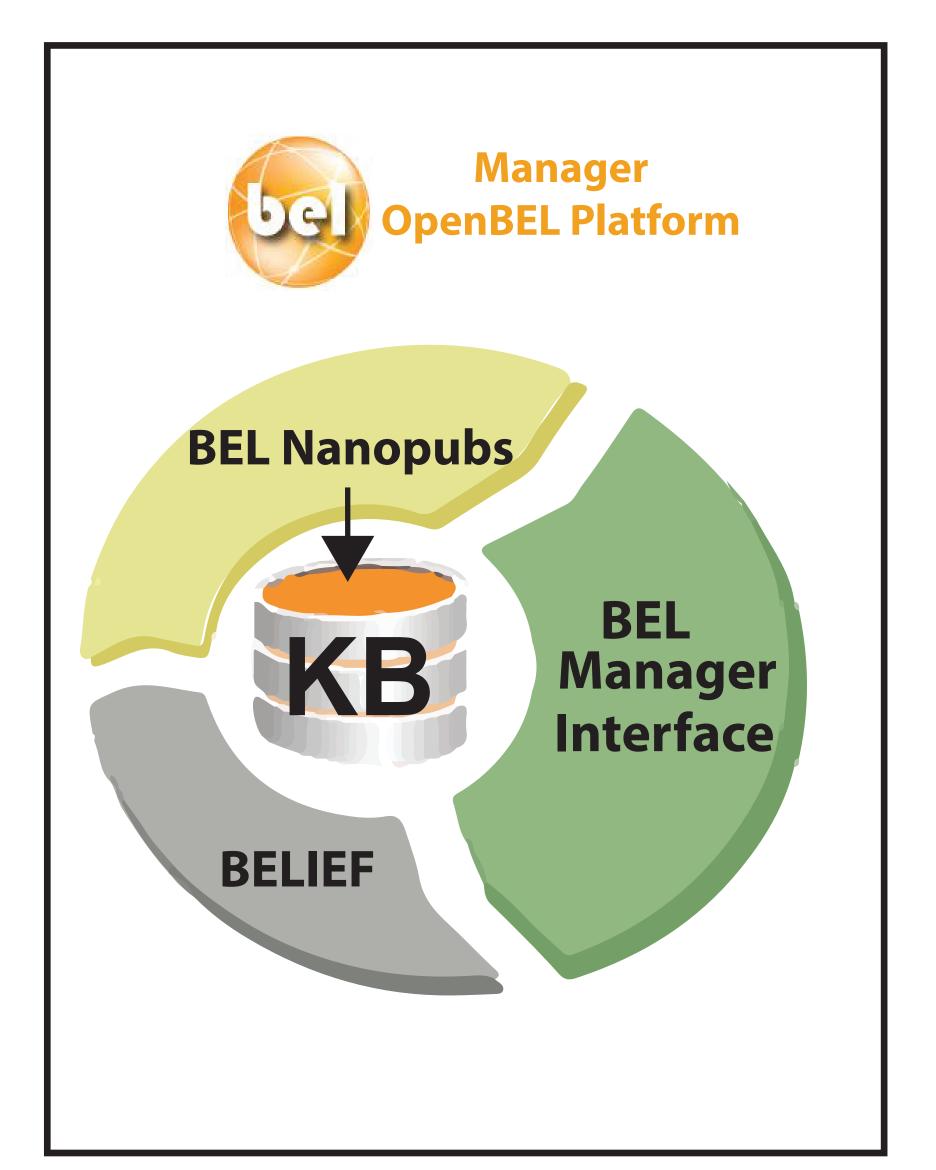


This schema describes the creation of structured knowledge base from the scientific literature. The workflow is initiated with the selection of scientific articles. After articles submission to the text mining pipeline, articles are processed, biological entities and relationships are extracted and assembled (1). The derived causal relationships are then compiled into a mechanistic network model (2). The network model describes molecular and cellular interactions that take place in vulnerable lesions accompanied with contextual information about the experiments (2). The systems toxicology approach developed by Philip Morris International (PMI) employs biological network models, transcriptomics data, and state of the art algorithms to quantify the impact of exposure on targeted biological processes (3-5).

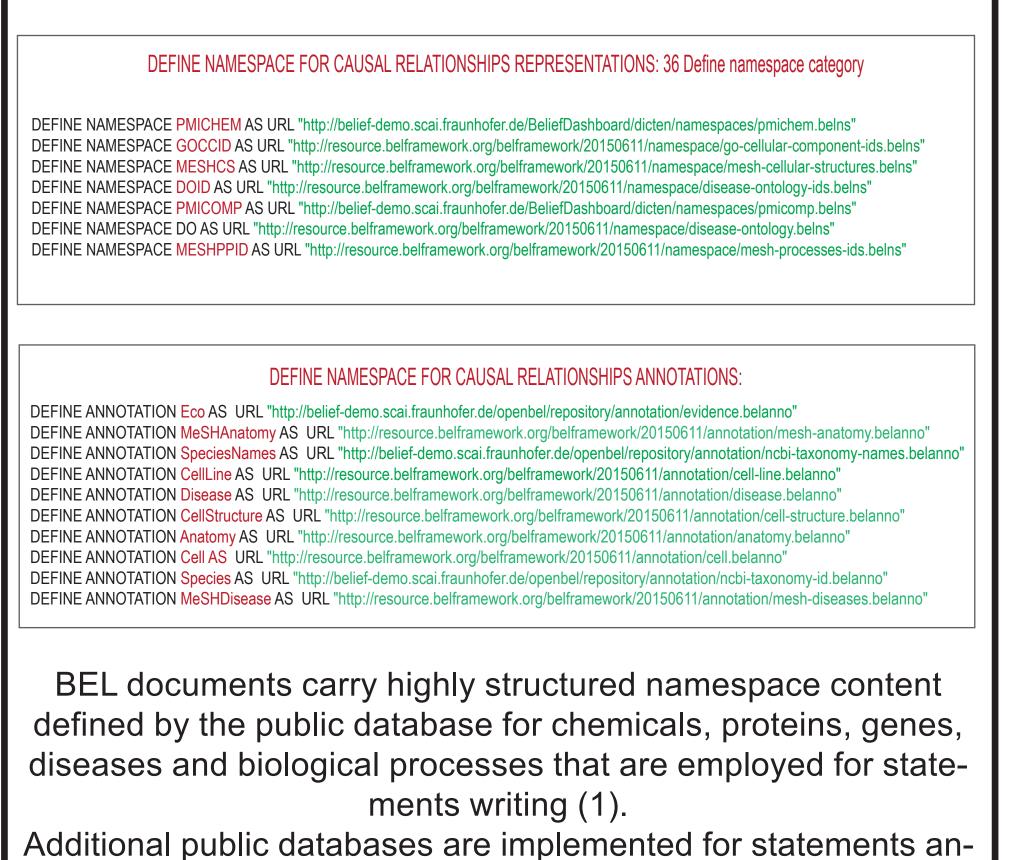
act(p(HGNC:MMP)) -> bp(GOBP:"Plaque Rupture",

BEL Manager: Knowledge Repository

Knowledge Repository

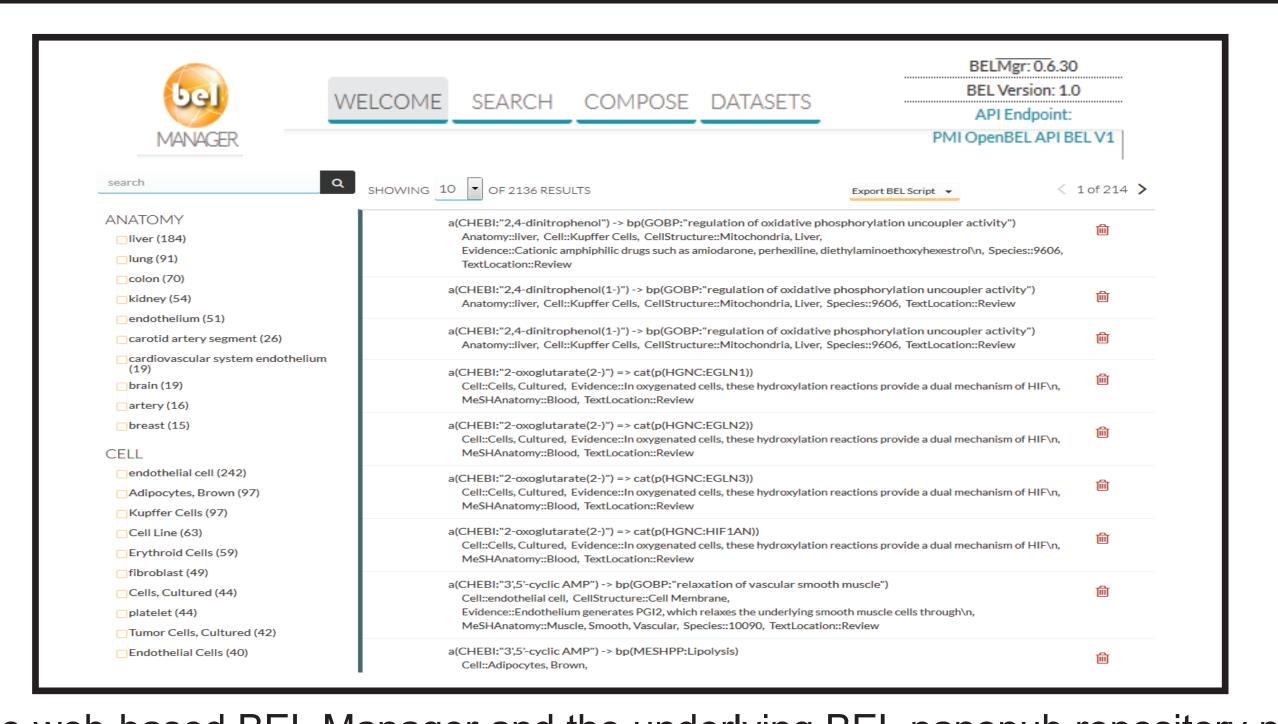


BEL Nanopubs



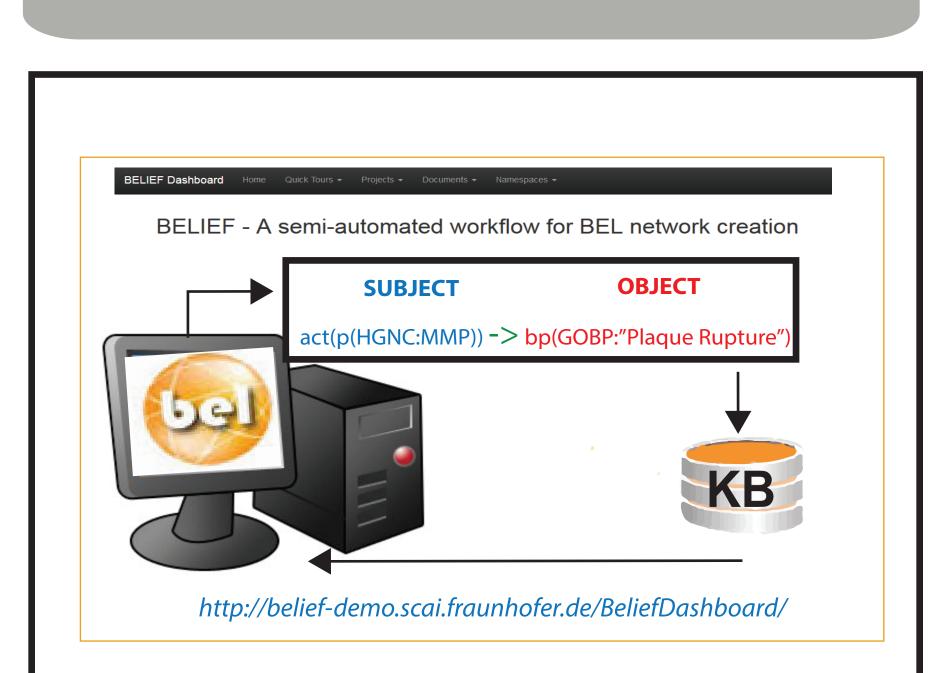
notations such as MESHAnatomy, MESHDisease, Eco...

BEL Manager Interface



The web-based BEL Manager and the underlying BEL nanopub repository provide access and allow storage of the curated relationships to the community. The BEL Manager enables the users to query and search the repository to export specific BEL contents for network refinement and reasoning.

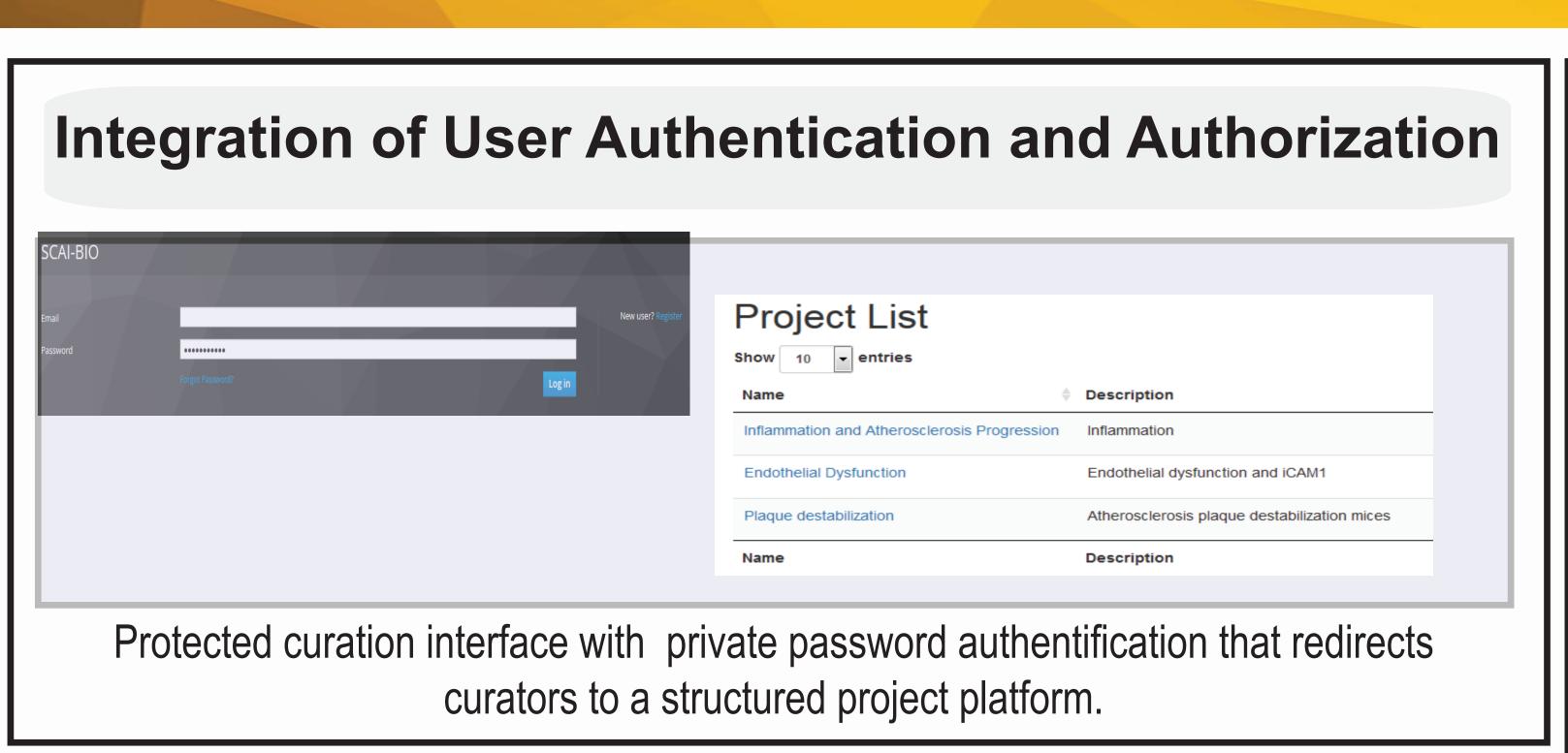
BELIEF



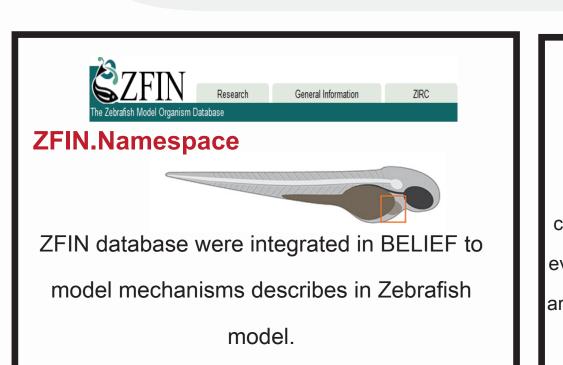
BELIEF (BEL Information Extraction Workflow)
has recently been upgraded via the integration
of the OpenBEL API enabling direct linkage to
the OpenBEL

platform and allowing for direct upload of curated documents into the BEL knowledge base (KB).

BELIEF Dashboard and Functions Improvements

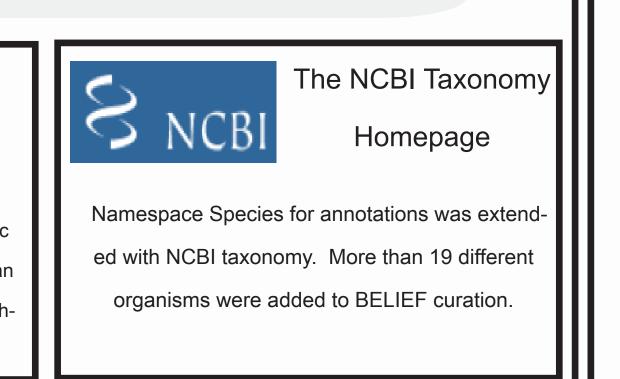


Expansion of Biomedical Vocabulary



THE EVIDENCE & CONCLUSION ONTOLOGY

The Evidence and Conclusion Ontology (ECO) is a controlled vocabulary that describes types of scientific evidence within the area of biological research that can arise from laboratory experiments, computational methods, manual literature curation, and other means.



Full text view area Authorised visualisation of all results part in the selected articles. Surrounding sentences in the full-text view could be selected and add as evidence.

Statements Annotations

Flexible annotation area allowed statement annotation at the evidence level.

Statement based statements:

Export

Exp

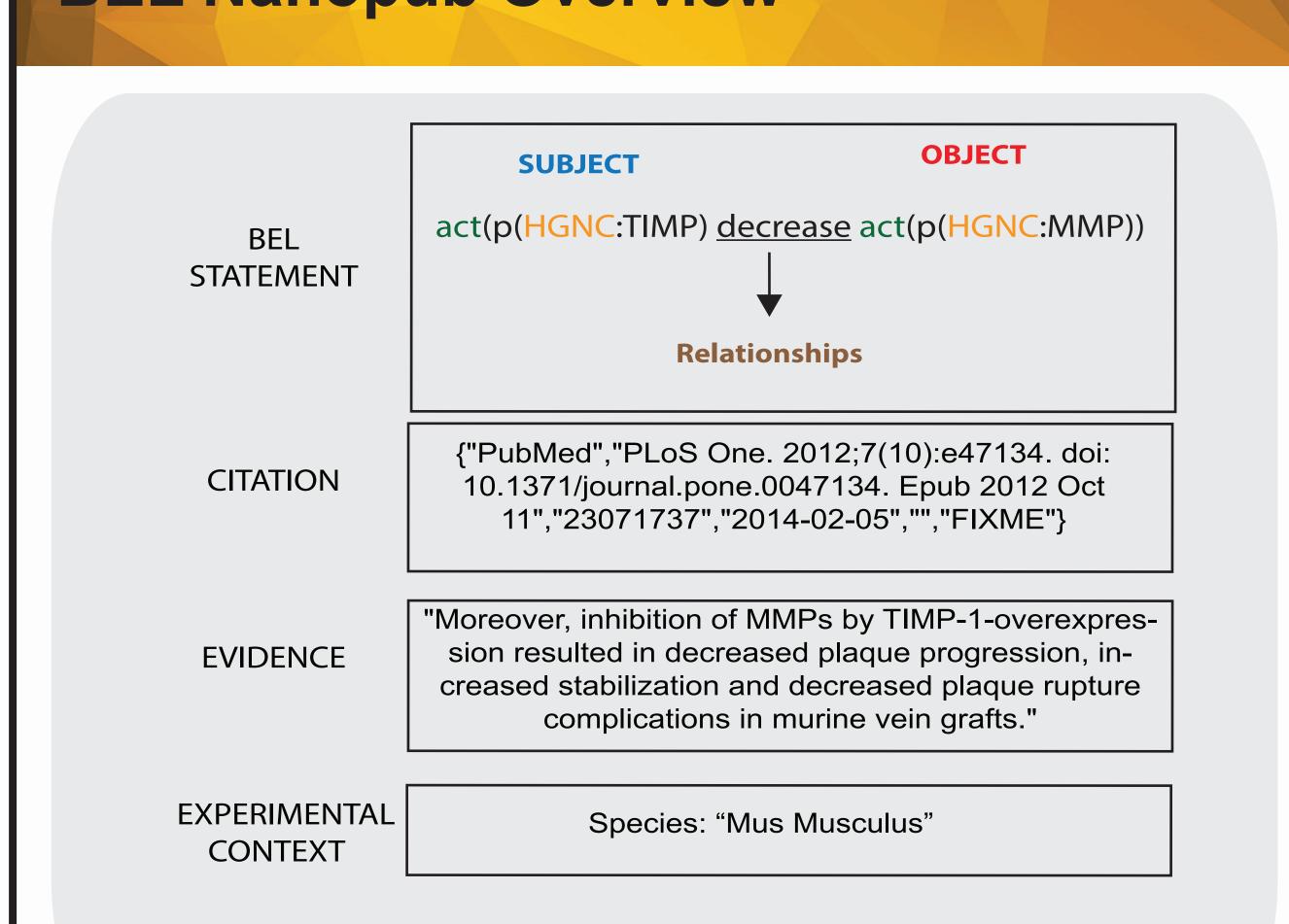
Search Namespace

Search concept tool box integrated into the BELIEF Dashboard, combine the curation process with searching activities on the dashboard.

Search namespaces

type e.g. CDK1

BEL Nanopub Overview



BEL NANOPUB

The three crucial elements of a BEL nanopub are the BEL statement showing the knowledge statement in a triple and controlled terminology, as well as the citation information and actual evidence sentence.

Experiment context is an additional field to simplify the triple assembly into biological network models (6).

Conclusions

Recent enhancements to BELIEF offer the opportunity to directly link the text mining tool and the curation process to a structured BEL Knowledge Base. A BEL nanopub represents the building block of the BEL KB. BEL Nanopub is the smallest unit of curation for BEL. We present here the integration of the concept of KB with the BELIEF text mining tool. The KB repository could be used to create, store, search and extract BEL nanopubs for network modelling or network refinement. Additional developments integrated into the BELIEF tool offer the opportunity to extend the mechanistic insight in KB and to maintain a deeply annotated contents.

References

Szostak J, Ansari S, Madan S, Fluck J, Talikka M, Iskandar A et al. Construction of biological networks from unstructured information based on a semi-automated curation workflow. Database: the journal of biological databases and curation 2015;2015:bav057. doi:10.1093/database/bav057.
 Szostak, J., Martin, F., Talikka, M., Peitsch, M.C., and Hoeng, J. (2016). Semi-Automated Curation Allows Causal Network Model Building for the Quantification of Age-Dependent Plaque Progression in ApoE-/- Mouse. Gene regulation and systems biology 10, 95-103.

(3) Hoeng J, Deehan R, Pratt D, Martin F, Sewer A, Thomson TM et al. A network-based approach to quantifying the impact of biologically active substances. Drug discovery today. 2012;17(9-10):413-8. doi:10.1016/j.drudis.2011.11.008.

(4) Iskandar AR, Xiang Y, Frentzel S, Talikka M, Leroy P, Kuehn D, et al. Impact assessment of cigarette smoke exposure on organotypic bronchial epithelial tissue cultures: a comparison of mono-culture and co-culture model containing fibroblasts. Toxicological Sciences. 2015:kfv122.

(5) Chindelevitch L, Ziemek D, Enayetallah A, Randhawa R, Sidders B, Brockel C et al. Causal reasoning on biological networks: interpreting transcriptional changes. Bioinformatics. 2012;28(8):1114-21. doi:10.1093/bioinformatics/bts090.

(6) Fluck, J., Madan, S., Ansari, S., Kodamullil, A.T., Karki, R., Rastegar-Mojarad, M., Catlett, N.L., Hayes, W., Szostak, J., Hoeng, J., et al. (2016). Training and evaluation corpora for the extraction of causal relationships encoded in biological expression language (BEL). Database: the journal of biological databases and curation 2016.