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

Justyna Szostak<sup>1</sup>, Sumit Madan<sup>2</sup>, William Hayes<sup>3</sup>, Jens Doerpinghaus<sup>2</sup>, Juliane Fluck<sup>2</sup>, Marja Talikka<sup>1</sup>, Manuel C. Peitsch<sup>1</sup>, Julia Hoeng<sup>1</sup>

<sup>1</sup> Philip Morris International R&D, Philip Morris Products S.A., Quai Jeanrenaud 5, 2000 Neuchâtel, Switzerland (Part of Philip Morris International group of companies). <sup>2</sup> Fraunhofer Institute for Algorithms and Scientific Computing, Schloss Birlinghoven, 53754 Sankt Augustin, Germany

<sup>3</sup>Applied 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.

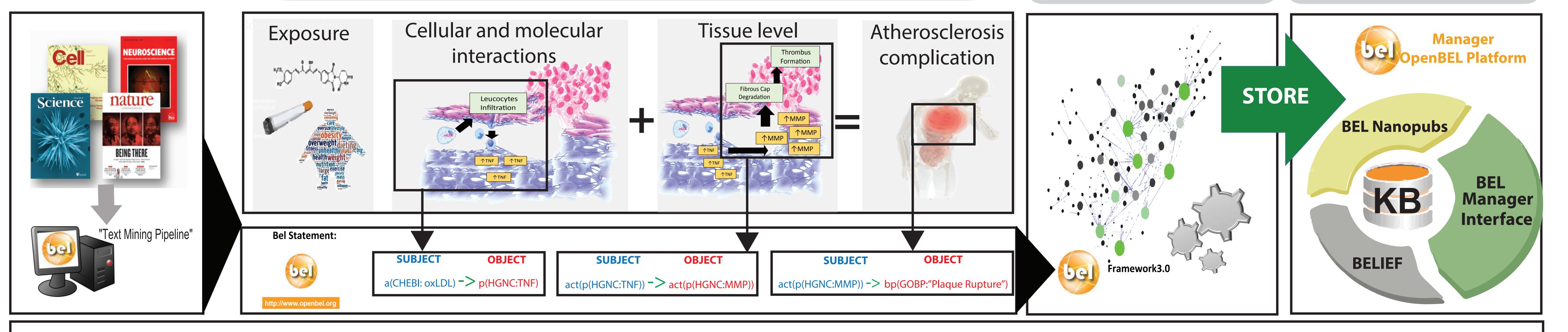
### **The BEL Information Extraction WorkFlow (BELIEF) Usecase**

Scientific Literature

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

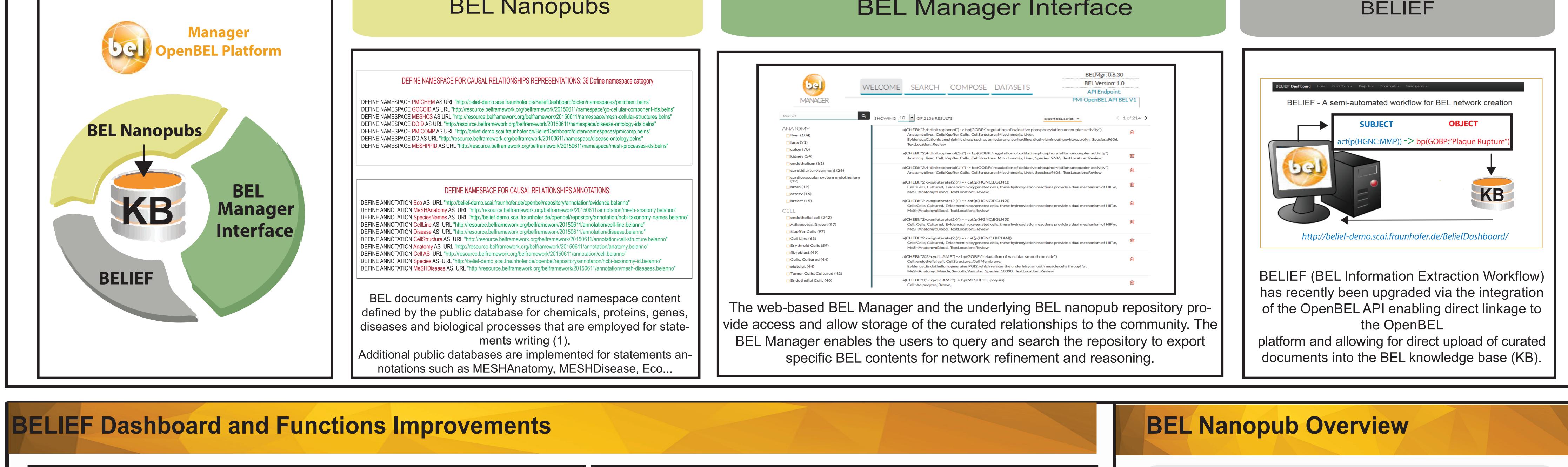
Structured Knowledge

**Knowledge Repository** 



This schema describes the creation of structured knowledge base from 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 bilogical processes (3-5).

# **BEL Manager : Knowledge Repository Knowledge Repository BEL Nanopubs BEL Manager Interface**



SCAI-BIO	
Email New user? Register	Project List Show 10 • entries
Forgot Password?	Name       Description         Inflammation and Atherosclerosis Progression       Inflammation

SCAI-BIO Email New user? Register	Project List Show 10 - entries	Full text view area of thrombosis The results showed that the area of thrombos was very fewer in control group (mean surface area was 0.9 ± 1.1 m m 2.) than that in model group (Table 1). Effect on Area of Thrombos The results showed that the area of thrombos was very fewer in control group (mean surface area was 0.9 ± 1.1 m m 2.) than that in model group (mean surface area was 78 ± 53 mm 2.) or in pravastatin group were lower than the structure of internal elastic men - brane was integral, tunica media was full of smooth muscle cells in control group. Whereas model group vessel showed inter-membrane incrassation, local ulcer in shoulder of the plaque, fibrin tissue hyperplasia and transparent denaturaliza - tion, and needle-like interspaces that resulted from cholesterol crystal dissolved by solvent. Foam cells filed with lipid vacu - ole were rich in the plaque, fibrin tissue hyperplasia in plaque, fibrin cap incrassation. Most shoulders of plaques were more integrated, plaque ulcers were smaller and medie-like interspaces in plaque were also diminished (Figs. 2a, b). Under electron microscopy, integrated endothelium, developed golgi complex, rough en -	Î STATEMENT	Relationships	
Forgot Password?	NameDescriptionInflammation and Atherosclerosis ProgressionInflammationEndothelial DysfunctionEndothelial dysfunction and iCAM1Plaque destabilizationAtherosclerosis plaque destabilization mices	results part in the selected         articles.         Surrounding sentences in the         Edit evidence selection	CITATION	{"PubMed","PLoS One. 2012;7(10):e47134. doi: 10.1371/journal.pone.0047134. Epub 2012 Oct 11","23071737","2014-02-05","","FIXME"}	
· · · · ·	NameDescriptionivate password authentification that redirectsuctured project platform.	full-text view could be selected and add as evidence.	EVIDENCE	"Moreover, inhibition of MMPs by TIMP-1-overexpres- sion resulted in decreased plaque progression, in- creased stabilization and decreased plaque rupture complications in murine vein grafts."	
Expansion of B	iomedical Vocabulary	Flexible annotation area allowed         statement annotation at the evi-	EXPERIMENTAL CONTEXT	Species: "Mus Musculus"	
The Zebrafish Model Organism Database	THE EVIDENCE & CONCLUSION ONTOLOGY The NCBI Taxonomy Homepage	dence level.		BEL NANOPUB	
ZFIN.Namespace       The Evidence and controlled vocabulation         ZFIN database were integrated in BELIEF to model mechanisms describes in Zebrafish       The Evidence and controlled vocabulation	A CONCLUSION NAMES A CONCLUSION		showing controlled terminolo Experiment of	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 evi- dence sentence. Experiment context is an additional field to simplify the triple assembly into biological network models (6).	

Integration of User Authentication and Authorization	Usability of Curation Dashboard	BEL STATEMENT	SUBJECT       OBJECT         act(p(HGNC:TIMP) decrease act(p(HGNC:MMP))
SCAI-BIO Email Password Forgot Password? Forgot Password?	Full text view area Full text view area Serum Lipid and Lipoprotein Levels of Pre- and Post- treatment for 4 weeks, serum TC, TG, and LDL-C levels in pravastatin group were lower than in model group (Table 1). Effect on Area of Thrombosis The results showed that the area of thrombus was very fewer in control group (mean surface area was 0.9 ± 1.1 m m 2.) than that in model group (mean surface area was 78 ± 53 mm 2.) or in pravastatin group were lower than in model group (mean surface area was 78 ± 53 mm 2.) or in pravastatin group were lower than the tructure of internal elastic media was fully decreased as compared with model group (Fig. 1). Effect on Morphology of Atheroscientic Aonta Light microscopy showed that the plaque, fibrin tissue hyperplasia and transparent denaturaliza - tion, and needle-like interspaces that resulted from the results of plaques were more integrated, plaque ulcers were smaller and there, foam cells and needle-like interspaces in plaque were also aliminished (Figs. 2a, b). Under electron microscopy, integrated endothelium, developed edgid complex, rough en - dopiamic reficultum, rich pinocytosis vesicle in plasma were visualized in ontorol group. In model group, endothelium was destructed obtiveliny, and endothelial cells could not be seen		Relationships
Name       Description         Inflammation and Atherosclerosis Progression       Inflammation         Endothelial Dysfunction       Endothelial dysfunction and iCAM1         Plaque destabilization       Atherosclerosis plaque destabilization mices	results part in the selected articles. Surrounding sentences in the Edit evidence selection	CITATION	{"PubMed","PLoS One. 2012;7(10):e47134. doi: 10.1371/journal.pone.0047134. Epub 2012 Oct 11","23071737","2014-02-05","","FIXME"}
Name         Description           Protected curation interface with private password authentification that redirects curators to a structured project platform.	full-text view could be selected and add as evidence.	EVIDENCE	"Moreover, inhibition of MMPs by TIMP-1-overexpres- sion resulted in decreased plaque progression, in- creased stabilization and decreased plaque rupture complications in murine vein grafts."
Expansion of Biomedical Vocabulary	Statements Annotations         Flexible annotation area allowed         statement annotation at the evi-	EXPERIMENTAL CONTEXT	Species: "Mus Musculus"
THE EVIDENCE & CONCLUSION ONTOLOGY The NCBI Taxonomy Homepage	dence level.	<b>BEL NANOPUB</b> 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 en- dence sentence. Experiment context is an additional field to simplify the triple assembly into biological network models (6).	
Verticative Model (regression beaution of the product of the prod	Search concept tool box integrated into the BELIEE Dash-		

## 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

(1) 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

(2) 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., 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.



**Biocuration 2017** Bringing knowledge to the people March 26-29, 2017 USA