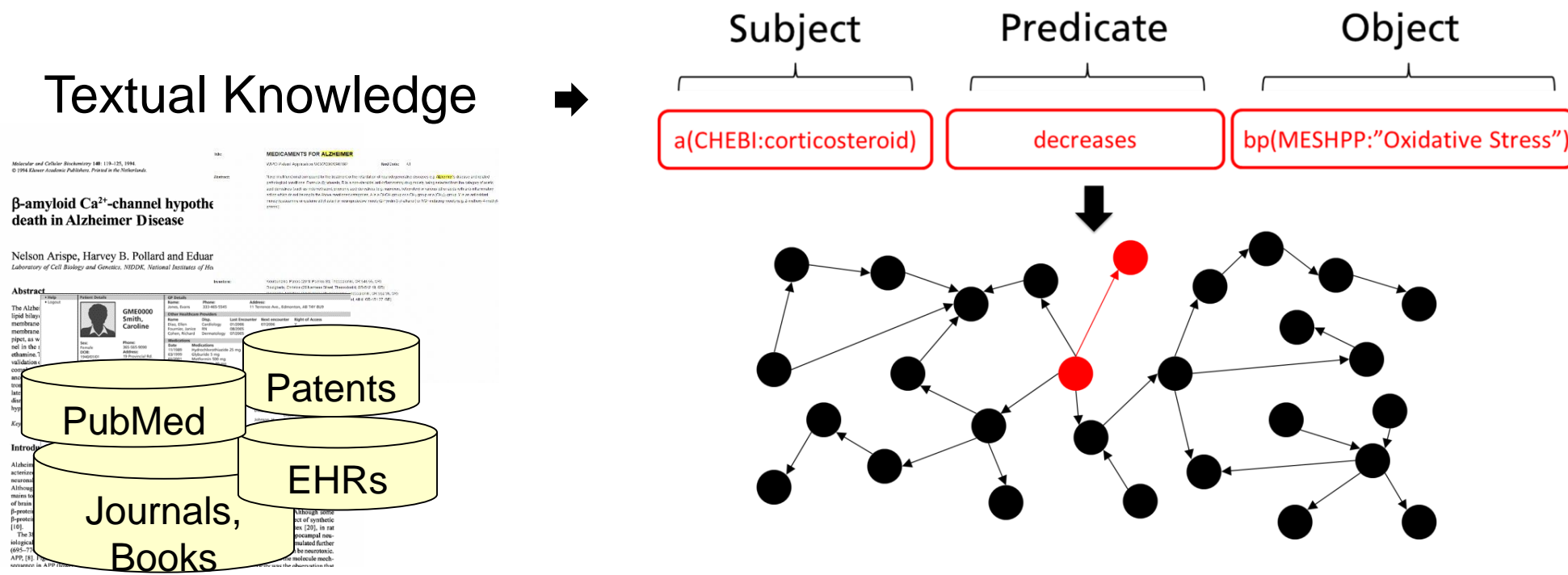# Bio-Medical Text Mining with Machine Learning

**Sumit Madan**

Department of Bioinformatics - Fraunhofer SCAI

# What is Bio-Medical Text Mining?

*Phosphorylation of glycogen synthase kinase 3 beta at Threonine, 668 increases the degradation of amyloid precursor protein*

Biological Expression Language (BEL)

p (HGNC:GSK3B, pmod (P,T,668)) -> deg (p (HGNC:APP))

BEL Functions

Namespace Identifiers

Relationship

Entity Definitions

# Biological Entity Classes

**Protein / Gene**
- HGNC
- ENTREZGENE
- UNIPROT
- MGI
- RGD
- …

**Small RNA**
- MIRBASE
- piRNABank
- …

Further Classes
- **Organism** (NCBITaxonomy)
- **Anatomy** (UBERON)
- **Phenotype** (HP)
- **Biological Process** (GO)
- **Cell** (CL)
- **Cell lines** (CLO)
- **Several clinical features**
- …

**Chemical**
- MESH
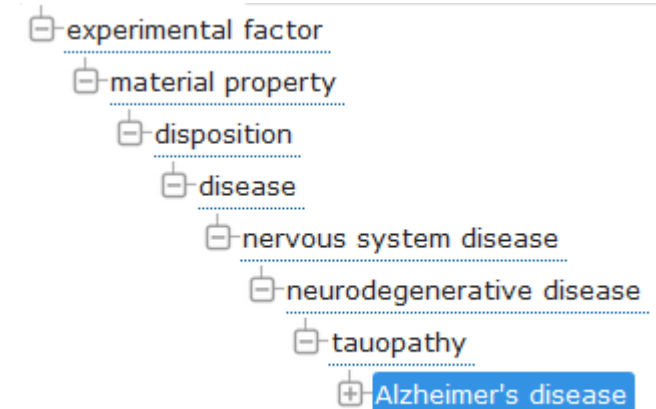- ChEBI
- ChEMBL
- …

**Disease**
- MESH
- DO
- ICD10
- …

**Bold: Entity Class**
Normal: Controlled Vocabulary

Fraunhofer
SCAI

# Several Spelling Variants (Synonyms) for One Single Concept



Over 28 Synonyms available just in Experimental Factor Ontology (EFO).

AD

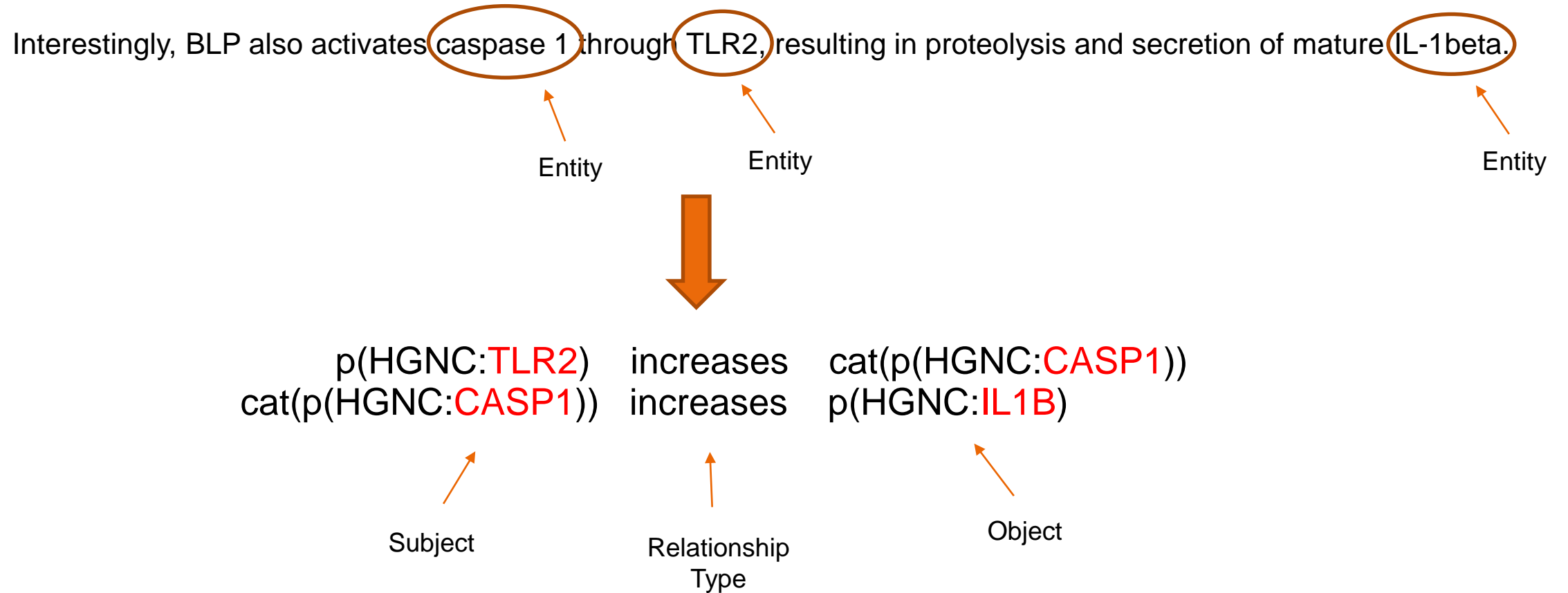Presenile Dementia

Alzheimer's Dementia

**Alzheimer's Disease** (EFO:0000249)

DAT

Alzheimer's disorder

Alzheimer's Disease (EFO:0000249)

experimental factor
  material property
    disposition
      disease
        nervous system disease
          neurodegenerative disease
            tauopathy
              Alzheimer's disease

Source: http://www.ebi.ac.uk/efo/EFO_0000249

Fraunhofer
SCAI

# Biological Relations (Example)

Interestingly, BLP also activates caspase 1 through TLR2, resulting in proteolysis and secretion of mature IL-1beta.

Entity                Entity                                                      Entity

p(HGNC:TLR2)              increases      cat(p(HGNC:CASP1))
cat(p(HGNC:CASP1))        increases      p(HGNC:IL1B)

Subject              Relationship                Object
                        Type

Fraunhofer SCAI

# Relation Extraction Workflow based on Convolutional Neural Network

Interestingly, BLP also activates caspase 1 through TLR2, resulting in proteolysis and secretion of mature IL-1beta (Source: PMID:10880445)
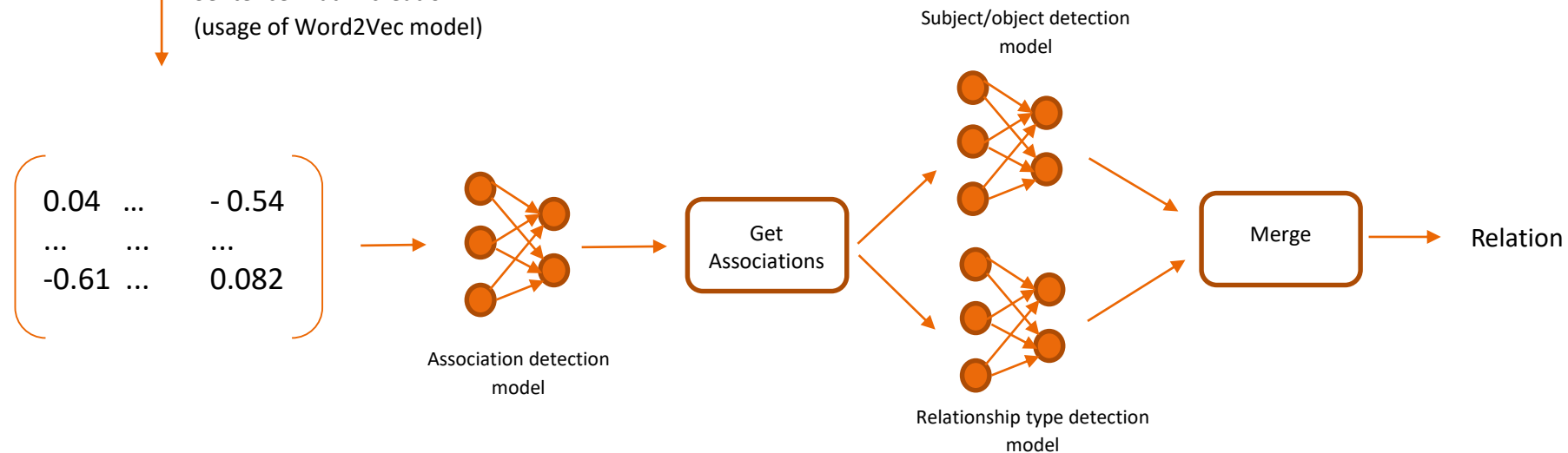
↓ Named Entity Recognition (NER)

Interestingly, BLP also activates caspase 1 through TLR2, resulting in proteolysis and secretion of mature IL-1beta

↓ Pre-processing

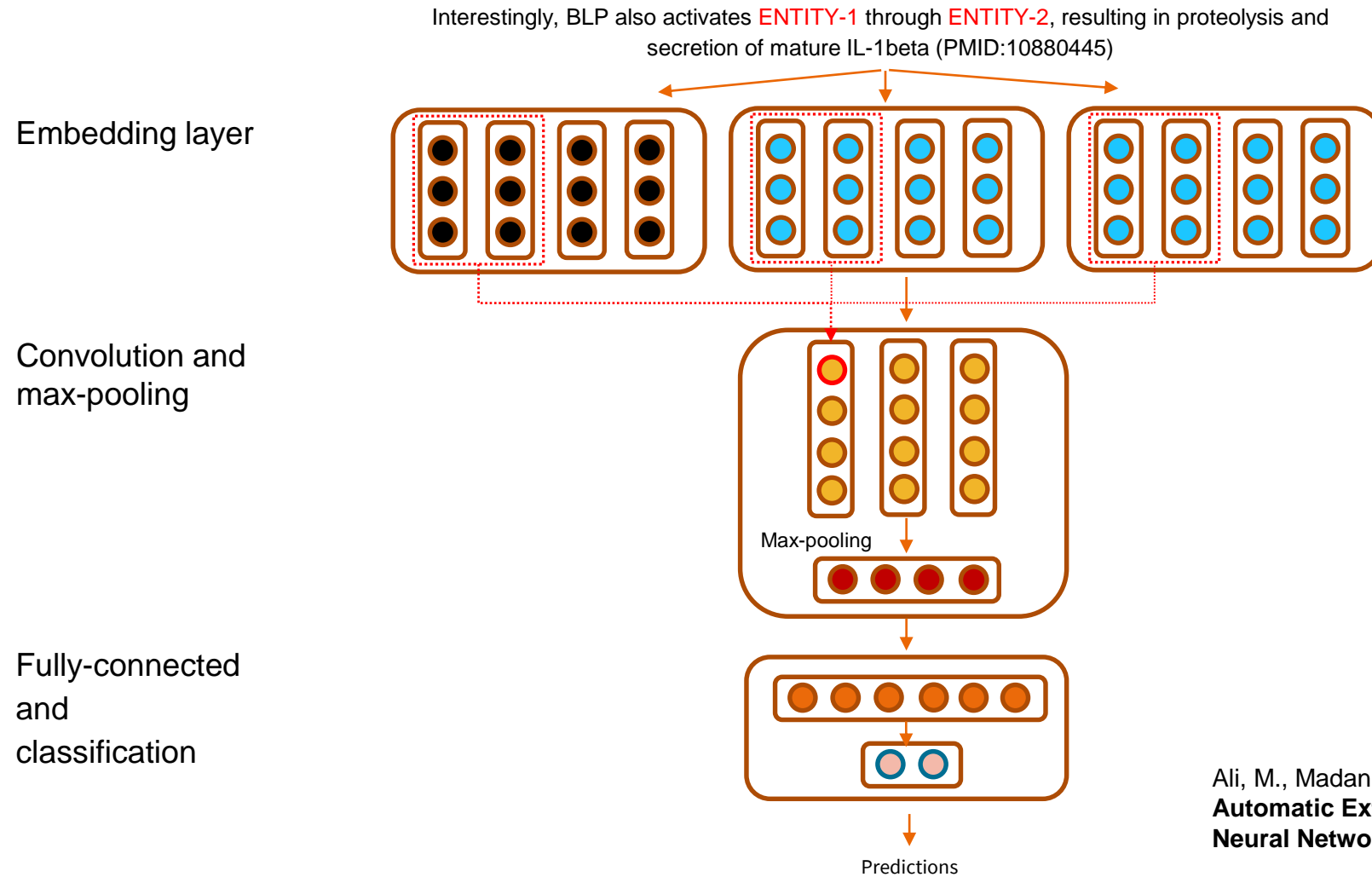Interestingly, BLP also activates ENTITY-1 through ENTITY-2, resulting in proteolysis and secretion of mature IL-1beta

↓ Sentence-matrix creation
(usage of Word2Vec model)



Ali, M., Madan, S., Fischer, A., et al. (2017) **Automatic Extraction of BEL-Statements based on Neural Networks**, *Proc. BioCreative VI Chall. Work.*

Fraunhofer

SCAI

# Multichannel Convolutional Neural Network (CNN) Architecture



Interestingly, BLP also activates ENTITY-1 through ENTITY-2, resulting in proteolysis and secretion of mature IL-1beta (PMID:10880445)

Embedding layer

Convolution and max-pooling

Max-pooling

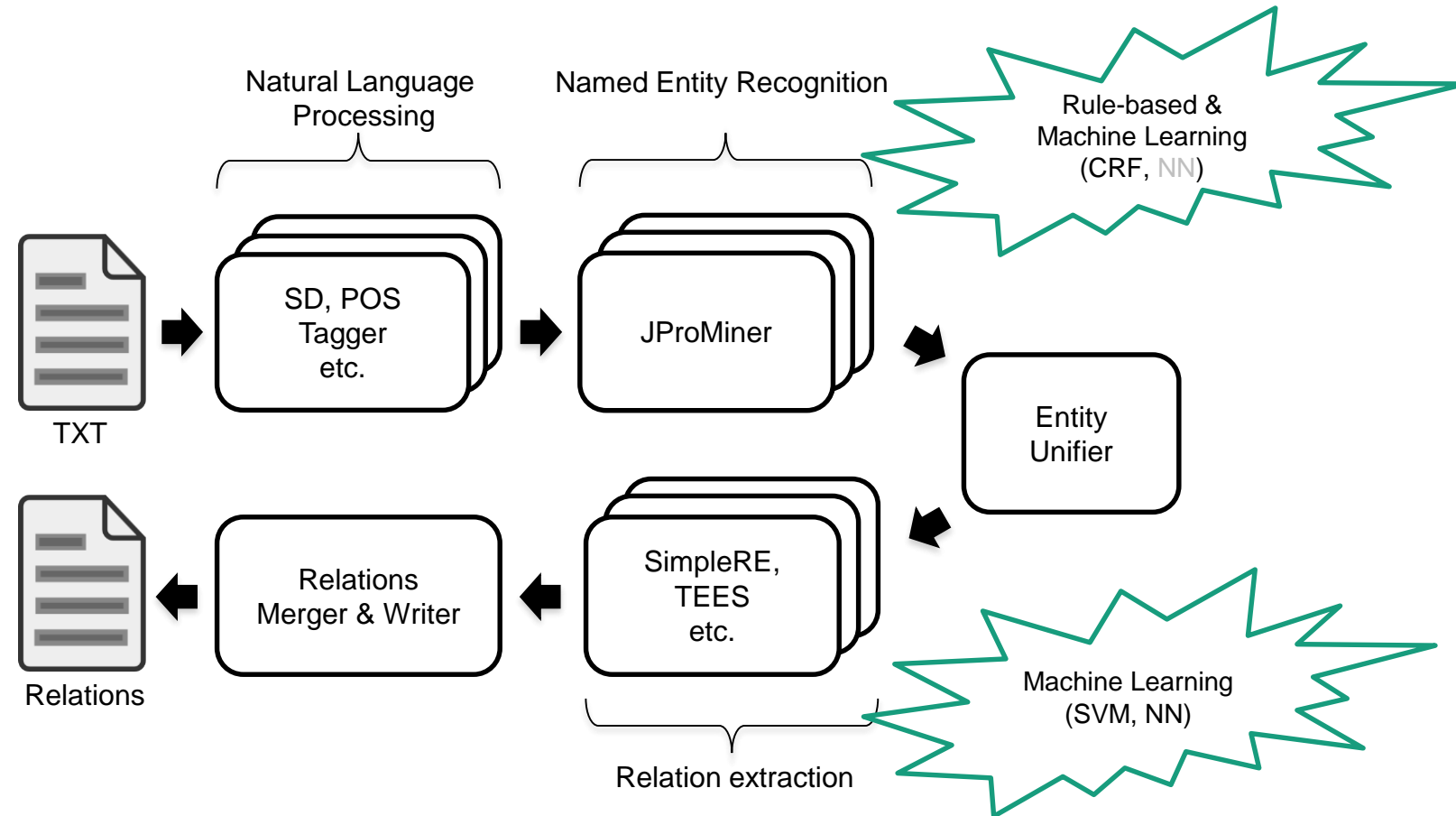Fully-connected and classification

Predictions

Ali, M., Madan, S., Fischer, A., et al. (2017) **Automatic Extraction of BEL-Statements based on Neural Networks**, *Proc. BioCreative VI Chall. Work.*

Fraunhofer

SCAI

# Larger Workflow for Relation Extraction (Project BELIEF)



Madan, S., Hodapp, S., Senger, P., et al. (2016) **The BEL information extraction workflow (BELIEF): evaluation in the BioCreative V BEL and IAT track**, *Database*, 2016, baw136.
http://belief.scai.fraunhofer.de/BeliefDashboard/

Fraunhofer
SCAI

# BELIEF Dashboard – Functionalities of Web-based Curation Interface

Curate document The novel KMO inhibitor CHDI-340246 leads to a restoration of electrophysiological alterations in mouse models of Huntington's disease.

← Return to document list | Go to statement centric view ➡

The novel KMO inhibitor CHDI-340246 leads to a restoration of electrophysiological alterations in mouse models of Huntington's ...
kynurenine (Kyn) pathway has been associated with the progression of Huntington's disease (HD). In particular, elevated levels of ...
3-hydroxy kynurenine (3-OH-Kyn) and quinol...
these metabolites is controlled by the activity ...
determine the role of KMO in the phenotype ...
this compound, when administered orally to t...
the central nervous system. The administrati...
Kynurenic acid (KynA) levels in brain tissues...
in mouse models of HD, both acutely and afte...
dosing of a selective KMO inhibitor does not ...

**Manually added evidences:**

**Export curated document**

| Curation Status | Curate | Export BEL | Reprocess |
|---|---|---|---|
| Approved | 🗄☑ | ⊕ ➤ | ↻ |
| Open | 🗄☑ | ⊕ | ↻ |
| Open | --- | --- | ↻ |

Edit evidence selection

**Statement section**
**Statements and annotations editing area**

Export ▾

Eco | Enter annotation value

1 (id:793): path(MESHD:"Huntington Disease") -- p(HGNC:KMO)

(id:1924): SpeciesNames | Mus musculus

(id:1925): MeSHDisease | Huntington Disease

Eco | Enter annotation value

**Detected concepts**

**Detected concepts**

**KMO**
HGNC:KMO ⤢
ZFIN:kmo ⤢
MGI:Kmo ⤢
**mouse**
SpeciesNames:"Mus musculus" ⤢
**Huntington's disease**
MeSHDisease:"Huntington Disease...

**Document Curation View**
(visualizes text mining results and also integrates several semantic search tools)

**Curation Status**

Status: Open
Last changed:
Comment:

✏ ✔ 👁 ☑ ➤

Comment:

Search namespaces

type e.g. CDK1 | ❶

**Statements search**

Add concepts and synonyms

Pubmed Information

Pubmed Id: 27163548 | Update

er
CAI

# Outlook

- Usage of bidirectional recurrent neural networks (BiLSTM)


- Use more training data from different tasks (such as BioNLP, and also BioCreative)

- Create further models to predict biological functions

- Automate hyper-parameter optimization

- Train new optimized word2Vec models (use PubMed, PMC & ontologies)

- Build hybrid systems: machine learning + rules-based system (UIMA Ruta)

**Fraunhofer**

**SCAI**