



Biomedical Text Mining

Fabio Rinaldi

Dalle Molle Institute for Artificial Intelligence (IDSIA) Lugano, Switzerland, USI/SUPSI

nlp.idsia.ch

November 16, 2021

Outline

Introduction

The Biomedical Literature

Methods

Our Tools (“OntoGene / IDSIA”)

References

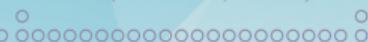
Introduction

The Biomedical Literature Methods



Our Tools ("OntoGene / IDSIA")

References



Topic

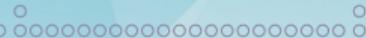
Introduction

The Biomedical Literature

Methods

Our Tools ("OntoGene / IDSIA")

References



What is text mining?

The ability to process text written in some human language (**unstructured data**), typically a large set of documents, interpret the meaning, and automatically extract concepts, as well as the relationships among those concepts, to directly answer questions of interest.

Unstructured data

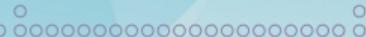
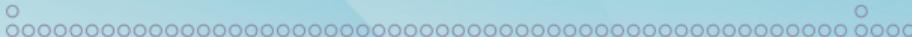
The university has 5600 students.
John's ID is number 1, he is 18 years old and already holds a B.Sc. degree.
David's ID is number 2, he is 31 years old and holds a Ph.D. degree. Robert's ID is number 3, he is 51 years old and also holds the same degree as David, a Ph.D. degree.

Semi-structured data

```
<University>
  <Student ID="1">
    <Name>John</Name>
    <Age>18</Age>
    <Degree>B.Sc.</Degree>
  </Student>
  <Student ID="2">
    <Name>David</Name>
    <Age>31</Age>
    <Degree>Ph.D. </Degree>
  </Student>
  ....
</University>
```

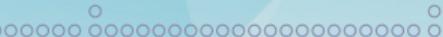
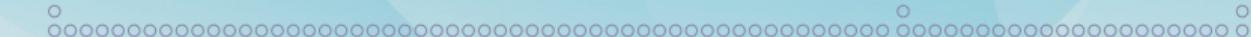
Structured data

ID	Name	Age	Degree
1	John	18	B.Sc.
2	David	31	Ph.D.
3	Robert	51	Ph.D.
4	Rick	26	M.Sc.
5	Michael	19	B.Sc.



What is unstructured data?

- Information that does not have a pre-defined data model (structure)
- Cannot be organized in a relational DB
- Typically text-heavy and usually contains free text in form of natural language
 - types: text, picture, videos (multimedia)
 - scientific literature, electronic health records, web pages, tweets and other social media text.
- These data sources/files may have internal structure but it is challenging to store information in row-column basis
- Estimates 80-90% of usable information in unstructured form
- Considerable growth and accumulation of unstructured data
- Often used as synonym of Big data (which is however both structured and unstructured)



Unstructured vs Structured Data



Structured Data

Often numbers or labels, stored in a structured framework of columns and rows relating to pre-set parameters.

ID CODES IN DATABASES

NUMERICAL DATA GOOGLE SHEETS

STAR RATINGS



Semi-unstructured Data

Loosely organized into categories using meta tags

EMAILS BY INBOX, SENT, DRAFT

TWEETS ORGANIZED BY HASHTAGS

FOLDERS ORGANIZED BY TOPIC



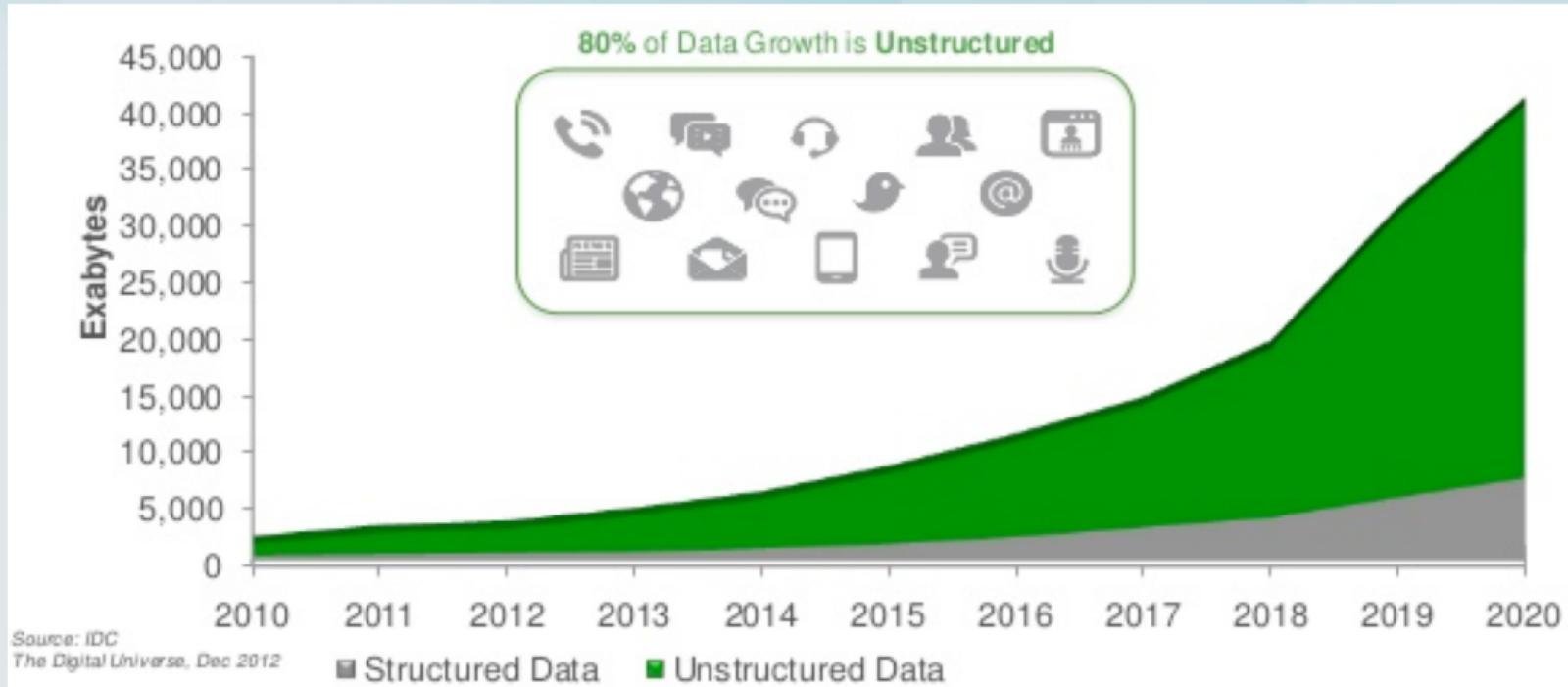
Unstructured Data

Text-heavy information that's not organized in a clearly defined framework or model.

MEDIA POSTS, EMAILS, ONLINE REVIEWS

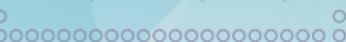
VIDEOS, IMAGES

SPEECH, SOUNDS

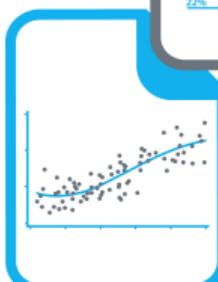
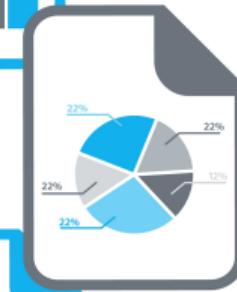
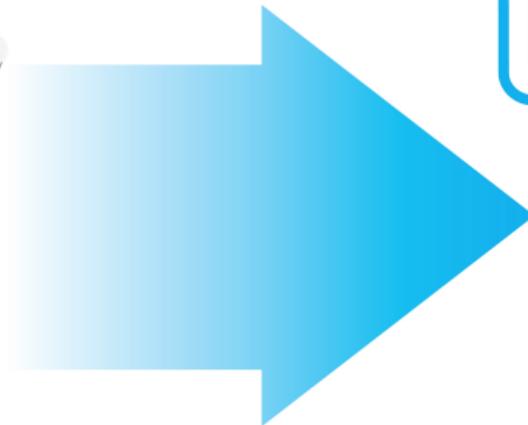
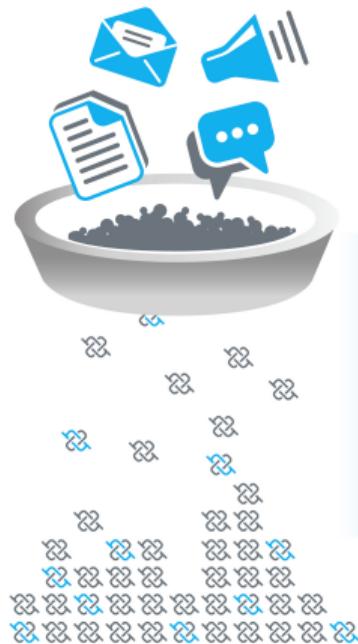


What is text mining?

- Data Mining needs structured data, usually in numerical form
- Text mining deals with unstructured data in form of natural language data (free text)
 - Goal: Retrieve what is hidden in text, present distilled knowledge in concise form
 - Process of discovery and extraction of knowledge from unstructured data
 - Extraction of non-trivial information or new information from natural language data collection
 - Applications integrating natural language processing components (such as document retrieval, classification and recognition of biological entities)

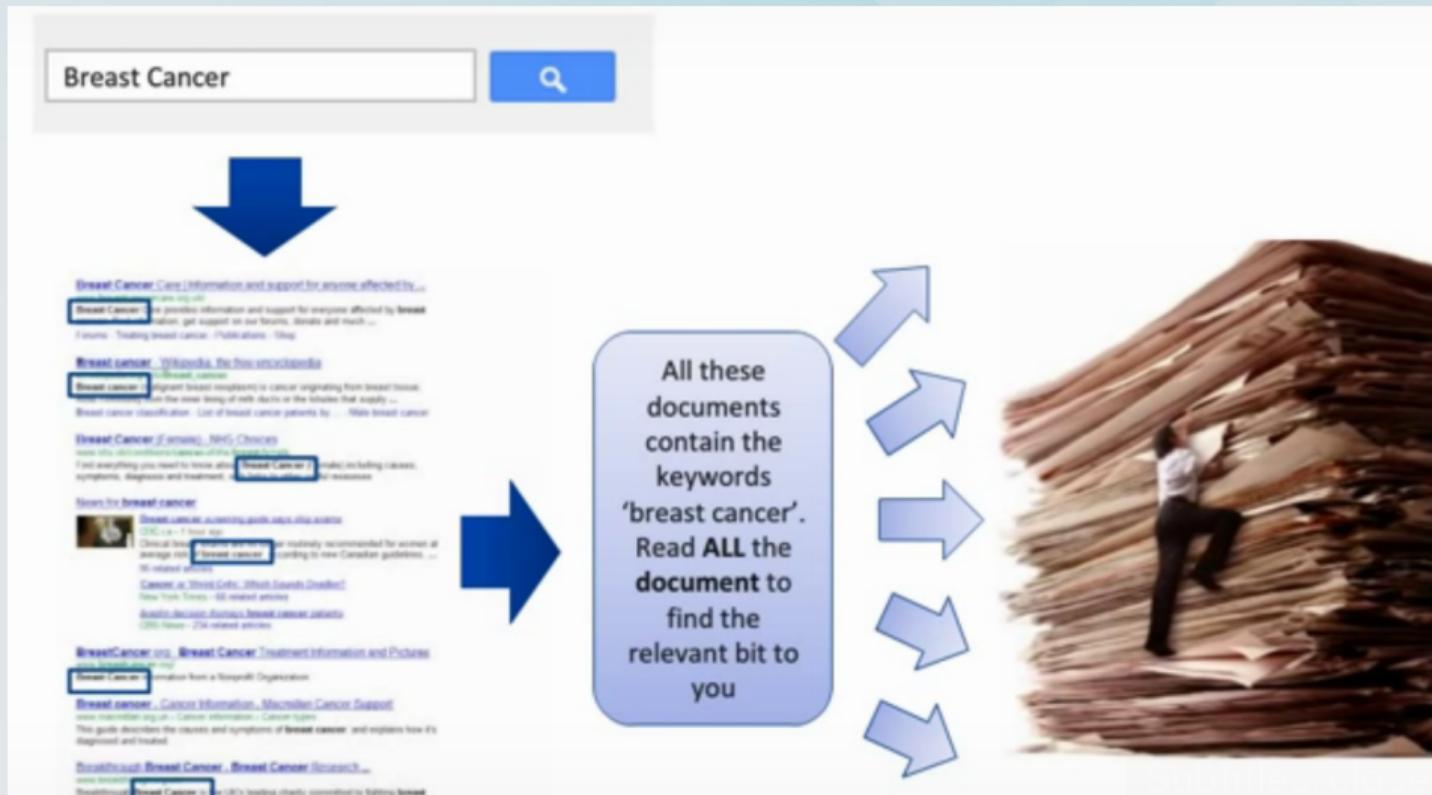


Text Mining process





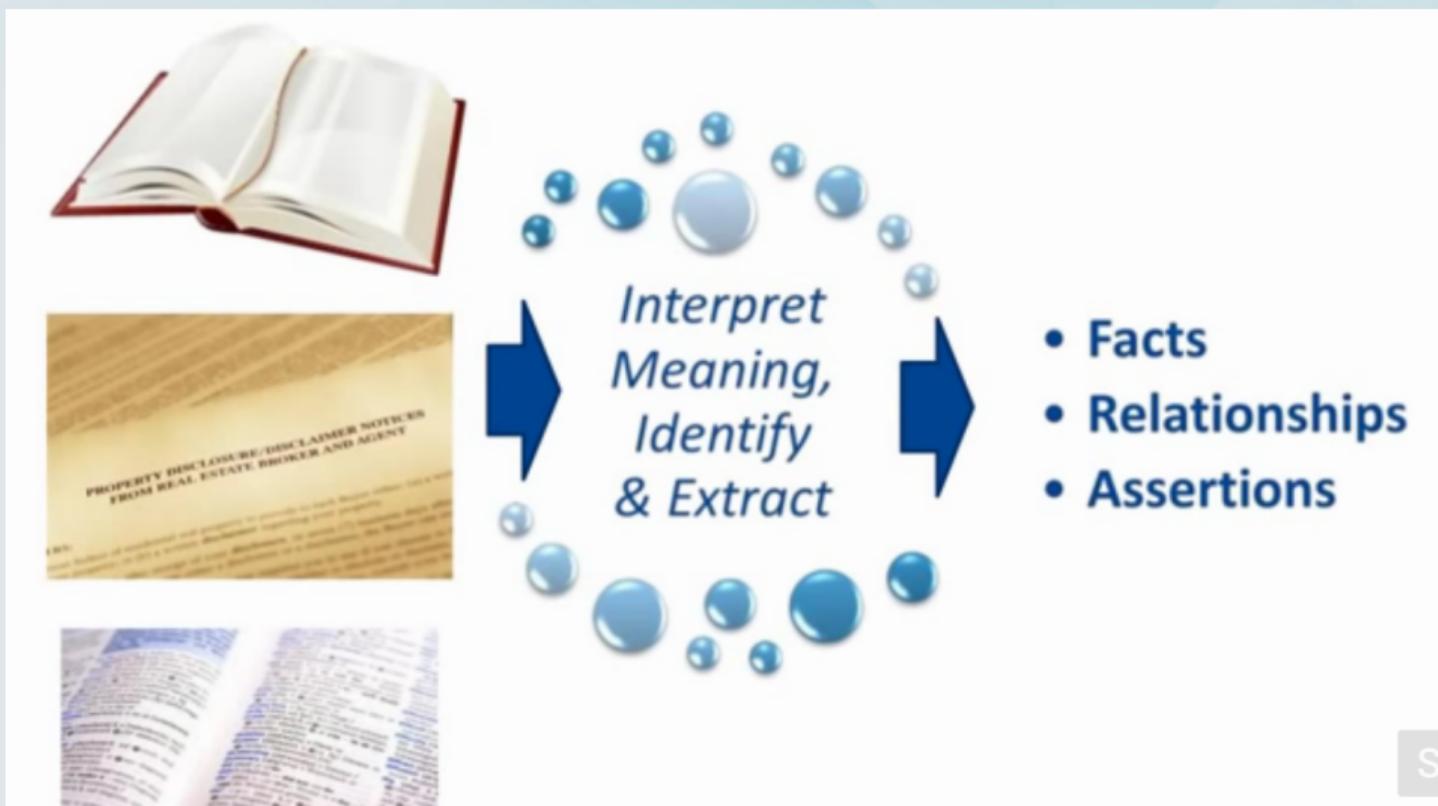
Conventional document search



Pros and cons

- Pros
 - speculative browsing for information
 - finding general information, e.g. address of the post office
 - Cons
 - large result sets, normally not enough time to read them all
 - noisy and irrelevant bits
 - need to consider variants of terms
 - e.g. *cancer, tumor, neoplasm*

What is Text Mining





Example

We find that p42mapk phosphorylates c-Myb on serine and threonine .

Purified recombinant p42 MAPK was found to phosphorylate Wee1 .

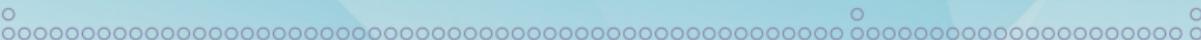
6:

Example

sentences

We find that p42mapk phosphorylates c-Myb on serine and threonine .

Purified recombinant p42 MAPK was found to phosphorylate Wee1 .



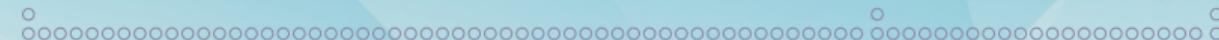
Example

sentences

noun groups
match *entities*

We find that p42mapk phosphorylates c-Myb or serine and threonine .

Purified recombinant p42 MAPK was found to phosphorylate Wee1 .



Example

sentences

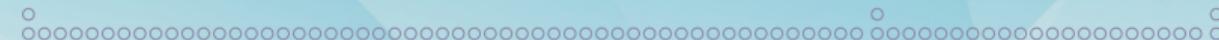
noun groups
match *entities*

verb groups
match *actions*

We find that p42mapk phosphorylates c-Myb or serine and threonine .

Purified recombinant p42 MAPK was found to phosphorylate Wee1 .

6:



Example

sentences

noun groups
match *entities*

verb groups
match *actions*

morphology -
different forms

We find that p42mapk phosphorylates c-Myb or serine and threonine .

Purified recombinant p42 MAPK was found to phosphorylate Wee1 .



Example

Genes affecting Breast Cancer? 

Hit

Estrogen-responsive finger protein as a new potential biomarker for breast cancer.

CONCLUSIONS: Our data suggest that Efp immunoreactivity is a significant prognostic factor in breast cancer patients. Moreover, Efp immunoreactivity was significantly correlated with poor prognosis of breast cancer patients, and multivariate analyses of disease-free survival and overall survival for 151 breast cancer patients showed that Efp immunoreactivity was the independent marker. Bcl-2 is a prognostic marker in breast cancer independently of the Nottingham Prognostic Index.

CONCLUSION: Bcl-2 is an independent predictor of breast cancer outcome and seems to be useful as a prognostic adjunct to the NPI, particularly in the first 5 years after diagnosis. Cyclin E as molecular marker in the management of breast cancer: a review.

Recent studies found cyclin E to be a promising prognostic indicator in breast cancer and examined its potential as a target for therapy.

NCOR1 mRNA is an independent prognostic factor for breast cancer.

We found risk for breast cancer was associated with the APOE genotype ($\chi^2 = 8.652$, $p = 0.013$).

Moreover, BRCA1 expression was elevated in HER4-positive human breast cancer specimens.

RESULTS: E-cadherin mRNA expression was lower in breast carcinoma ($P = 0.001$), whereas Snail expression was higher ($P = 0.003$).

Estrogen-responsive finger protein as a new potential biomarker for breast cancer

Efp immunoreactivity was significantly correlated with breast cancer patients

Bcl-2 is an independent predictor of breast cancer outcome

BRCA1 expression was elevated in HER4-positive human breast cancer

Cytoplasmic CLU expression

Se

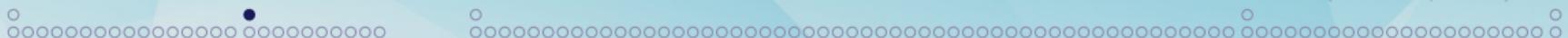
Introduction

The Biomedical Literature Methods



Our Tools ("OntoGene / IDSIA")

References



Topic

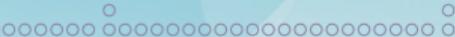
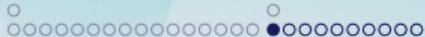
Introduction

The Biomedical Literature

Methods

Our Tools ("OntoGene / IDSIA")

References

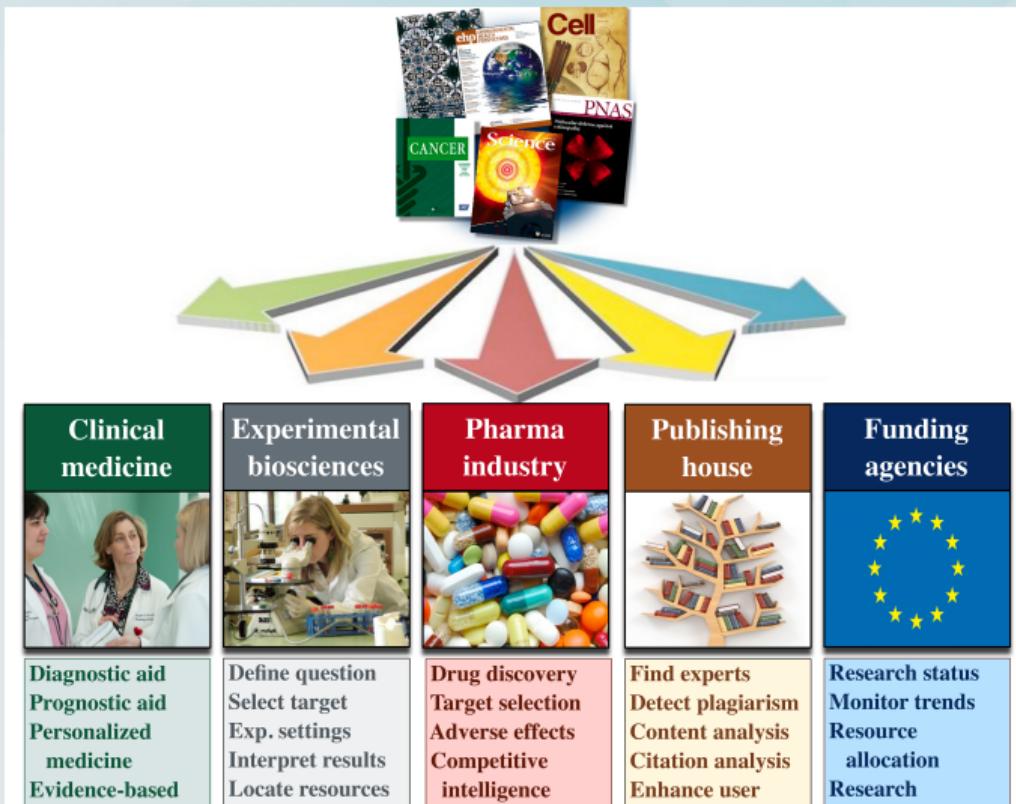


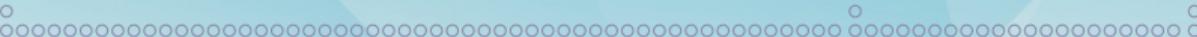
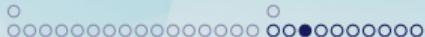
Biomedical literature

- Huge number of documents available
 - ~28 million references (Medline)
 - more than half of them have an abstract in English
 - two publications every minute
 - about 1.5 million full text (PMC)
- Aims of literature mining
 - locate and extract relevant biological information
 - integrate and cross-link it with other biological data and knowledge bases



Secondary usage of the literature





Major resource: Medline/PubMed

- One of the resources maintained by the National Library of Medicine
- Origin: Index Medicus, 1879-2004
- Going digital: MEDLARS, 1964-

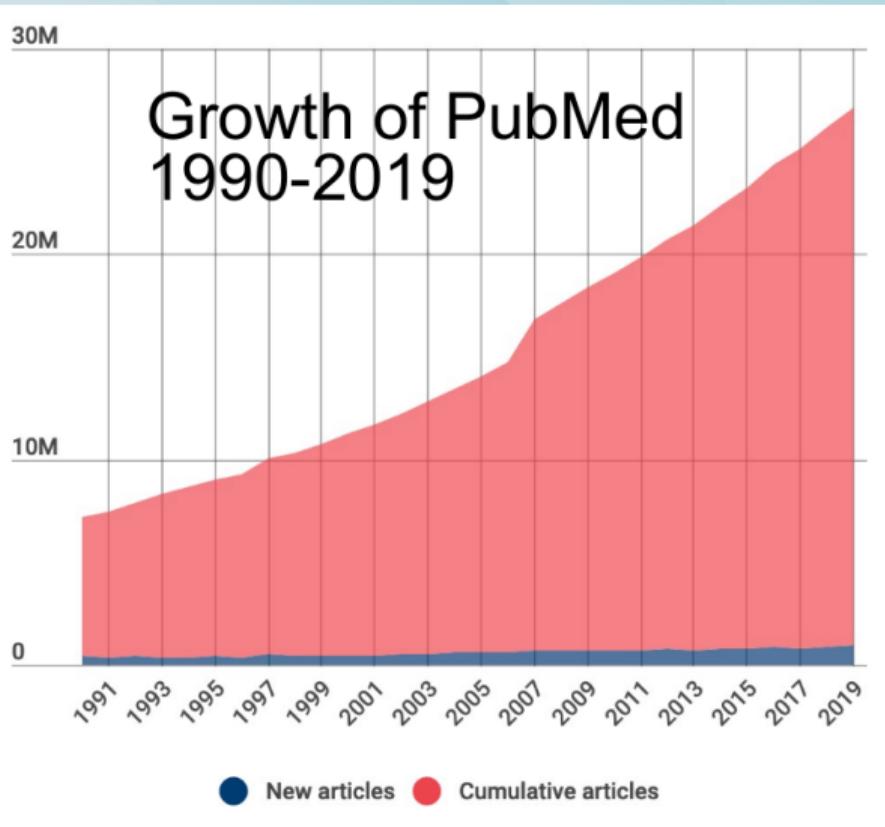
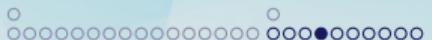
MEDLARS (Medical Literature Analysis and Retrieval System) is a computerised biomedical bibliographic retrieval system.

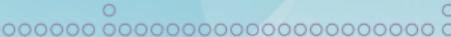
- Going online: MEDLINE, 1971-

In late 1971, an online version called MEDLINE ("MEDLARS Online") became available as a way to do online searching of MEDLARS from remote medical libraries. This early system covered 239 journals and boasted that it could support as many as 25 simultaneous online users.

- Going on the web: PubMed, 1997-

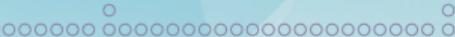
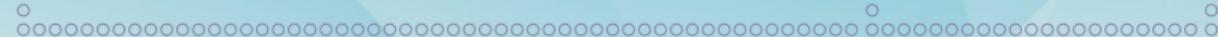
more than 28 million records from 5,639 selected publications covering biomedicine and health from 1950 to the present.





Open Access

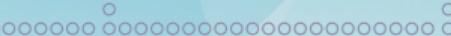
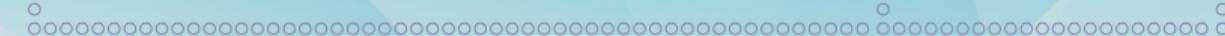
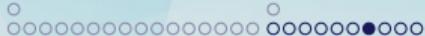
- 2003, Bethesda Statement on Open Access Publishing
(<http://www.earlham.edu/~peters/fos/bethesda.htm>)
- 2003, Berlin Declaration on Open Access to Knowledge in the Sciences and Humanities
<https://openaccess.mpg.de/Berlin-Declaration>
- 2004, PubMedCentral (PMC) <https://www.ncbi.nlm.nih.gov/pmc/>
 - 4.5 million articles
 - 2027 “Full participation” journals
 - ~4700 Partial participation journals



Open Access

- NIH requirement

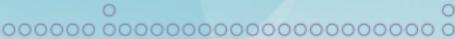
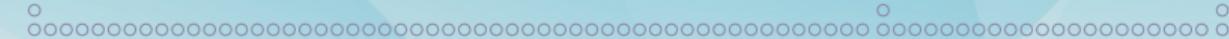
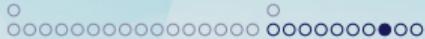
- "Policy on Enhancing Public Access to Archived Publications Resulting From NIH-Funded Research,"
- all publications generated by NIH-funded research (in whole or part) must be submitted to PMC
- from May 2, 2005
- NLM is digitizing earlier print issues of many of the journals already in PMC, extending the availability of full texts back to before the implementation of the 2005 policy.



NOT Open Access

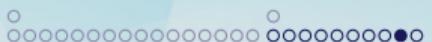
Unfortunately, free full-text access remains impossible for a large portion of scientific journals. In some fields, such as chemistry, even article abstracts are inaccessible for a large-scale analysis. The obvious outcome is that articles published in open-access journals have a better chance of being identified as relevant hits than others appearing in traditional "closed-access" journals. Electronic access to text obviously impacts all stages of text mining.

[Rzhetsky et al., 2009]

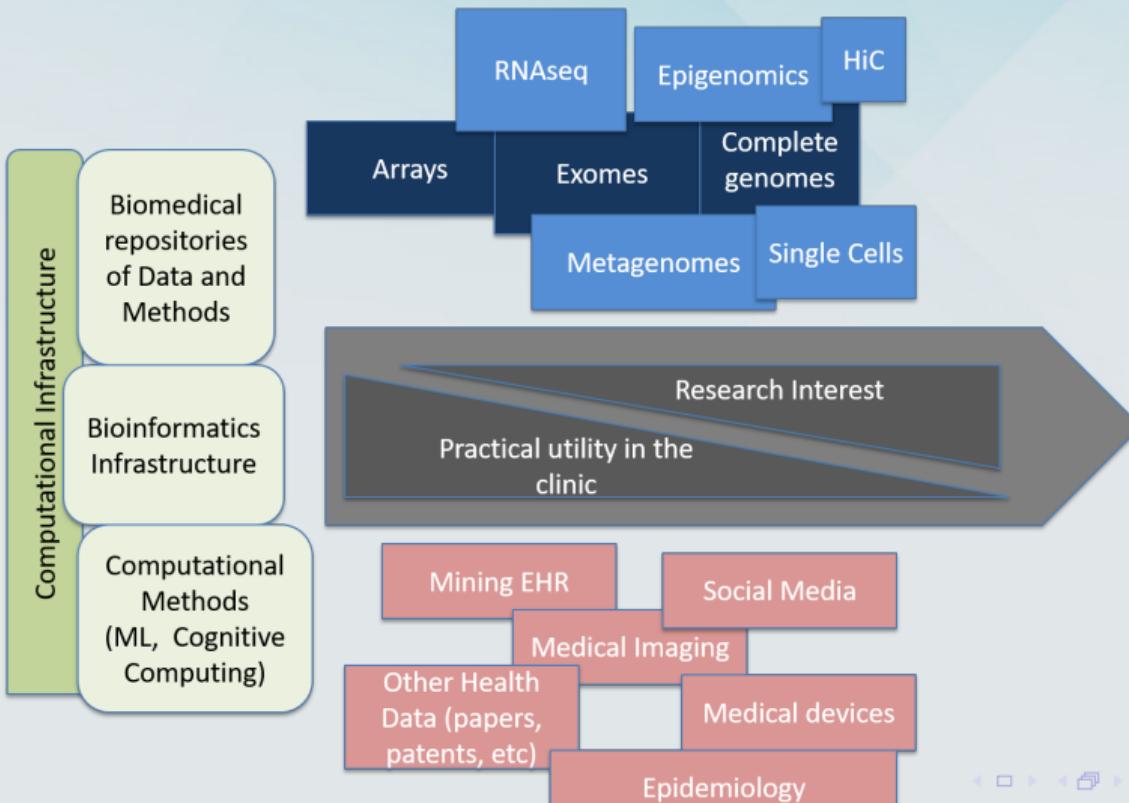


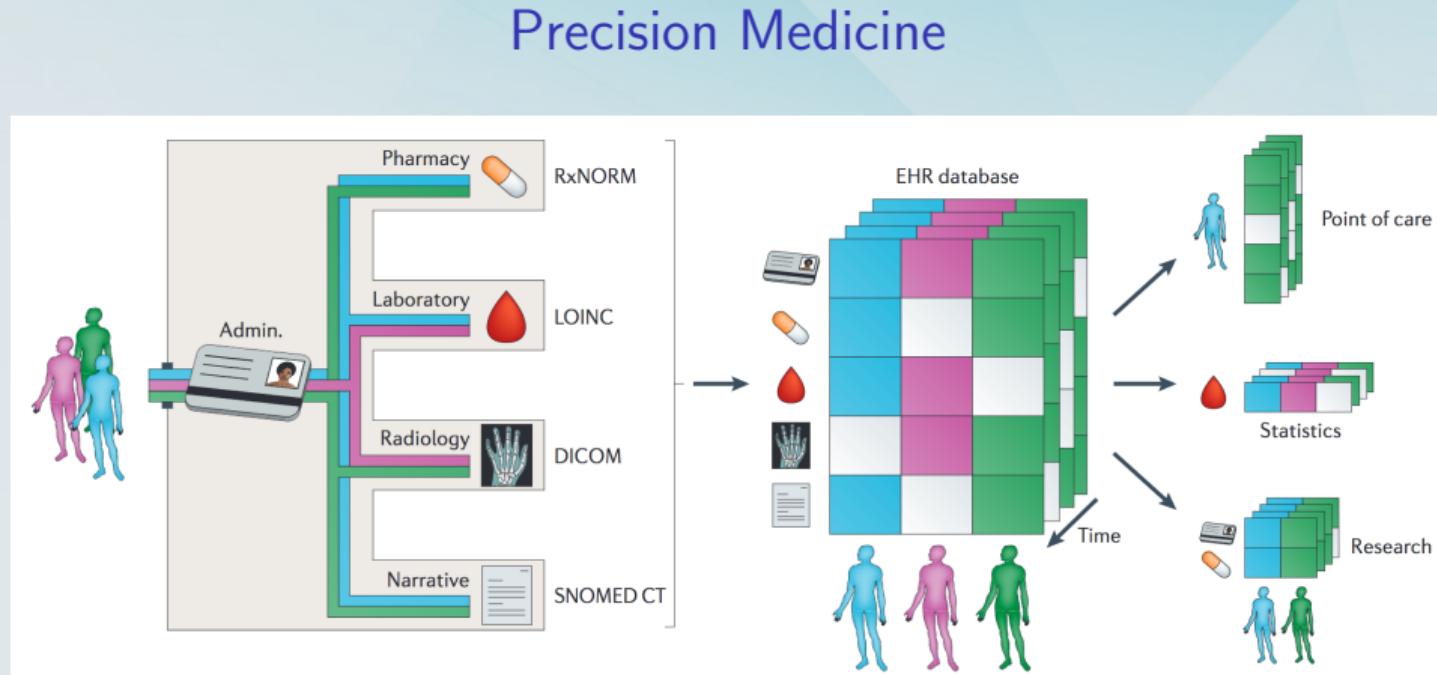
Why PMC?

PMC “makes it possible to integrate the literature with a variety of other information resources such as sequence databases and other factual databases that are available to scientists, clinicians and everyone else interested in the life sciences. The intentional and serendipitous discoveries that such links might foster excite us and stimulate us to move forward.”



Not only PubMed: Precision Medicine





Jensen, P. B., Jensen, L. J., & Brunak, S. (2012). Mining electronic health records: towards better research applications and clinical care. *Nature reviews. Genetics*, 13(6), 395.

Introduction

The Biomedical Literature

Methods



Our Tools ("OntoGene / IDSIA")

References

Topic

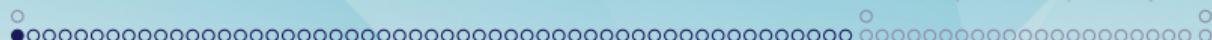
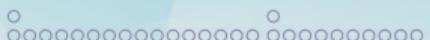
Introduction

The Biomedical Literature

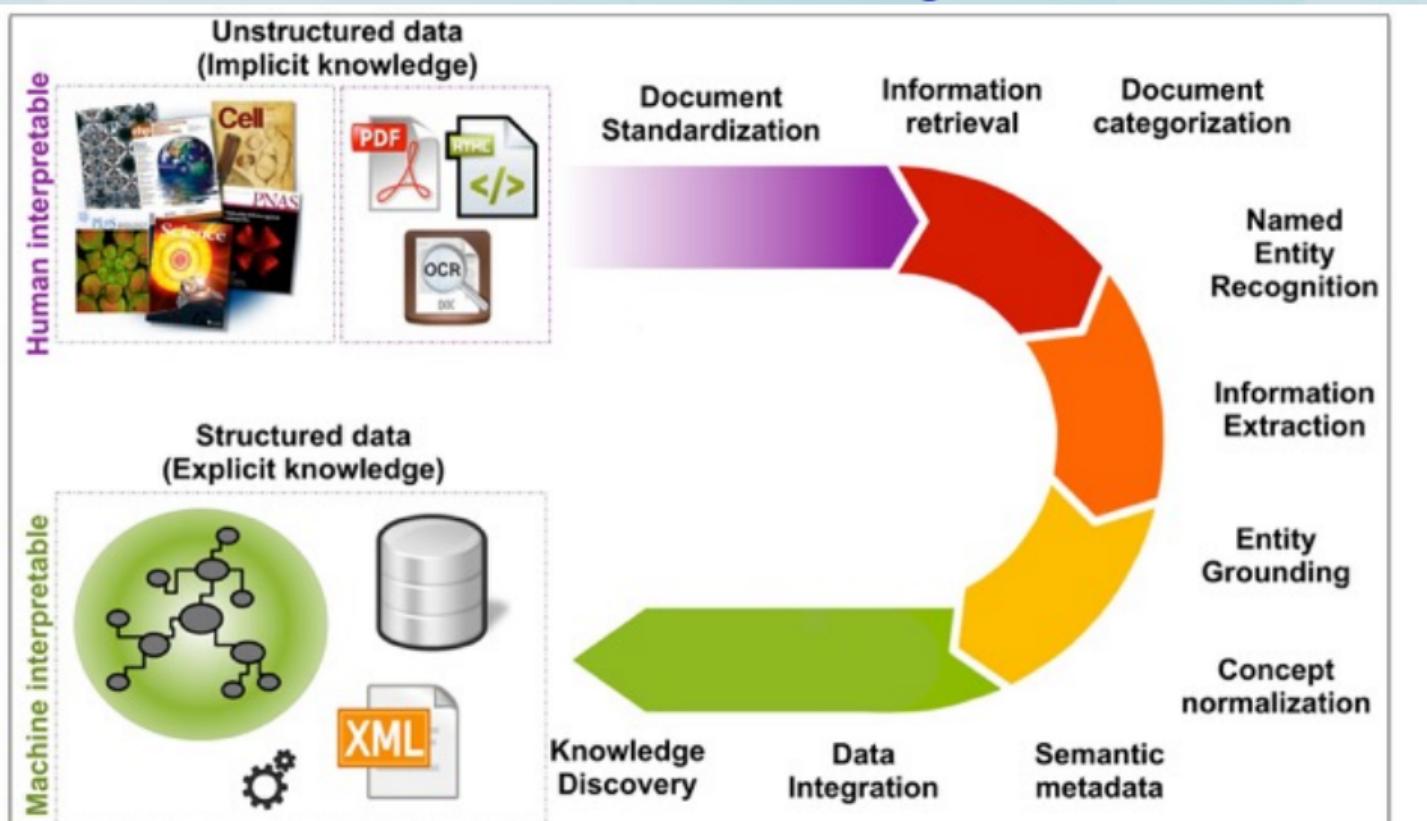
Methods

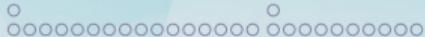
Our Tools ("OntoGene / IDSIA")

References



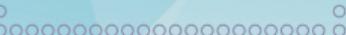
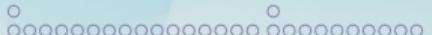
Canonical Text Mining flowchart





IR, IE, DM

- Information Retrieval (IR)
 - find relevant documents (unstructured)
 - example: PubMed search
- Information Extraction (IE)
 - find detailed information (semi-structured)
 - NER: Named Entity Recognition
 - Concept Recognition
 - Relation Extraction
- Data Mining (DM)
 - find associations, build networks, make predictions



PubMed as an IR engine

NCBI Resources How To My NCBI Sign In

PubMed.gov U.S. National Library of Medicine National Institutes of Health

Search: PubMed RSS Save search Limits Advanced search Help

gastrin Search Clear

Display Settings: Summary, 20 per page, Sorted by Recently Added Send to:

Results: 1 to 20 of 19878 << First < Prev Page 1 Next > Last >>

[The Thyroid and the Gut.](#)
1. Ebert EC.
J Clin Gastroenterol. 2010 Mar 25. [Epub ahead of print]
PMID: 20351569 [PubMed - as supplied by publisher]
[Related articles](#)

[Preclinical and clinical studies of peptide receptor radionuclide therapy.](#)
2. Pool SE, Krenning EP, Koning GA, van Eijck CH, Teunissen JJ, Kam B, Valkema R, Kwekkeboom DJ, de Jong M.
Semin Nucl Med. 2010 May;40(3):209-18.
PMID: 20350630 [PubMed - in process]
[Related articles](#)

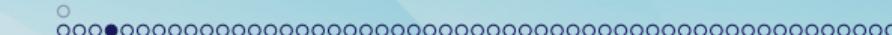
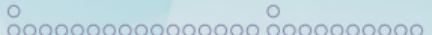
[Gastrin response to candidate messengers in intact conscious rats monitored by antrum microdialysis.](#)
3. Ericsson P, Häkanson R, Norlén P.
Regul Pept. 2010 Mar 24. [Epub ahead of print]
PMID: 20346991 [PubMed - as supplied by publisher]
[Related articles](#)

[Dyspeptic Symptom Development After Discontinuation of a Proton Pump Inhibitor: A Double-Blind Placebo-Controlled Trial.](#)
4. See more

Filter your results:
All (19878)
Review (1740)
Free Full Text (2831)
Manage Filters

Also try:
▶ gastrin releasing peptide receptor
▶ serum gastrin
▶ gastrin gastric
▶ gastrin level
▶ gastrin review

Titles with your search terms
▶ Combination therapy with glucagon-like peptide-1 and gastrin restore [Diabetes. 2008]
▶ Importance of gastrin in the pathogenesis and treatment of [World J Gastroenterol. 2009]
▶ Combination therapy with glucagon-like peptide-1 and gastrin i [Cell Transplant. 2008]



PubMed as an IR engine

The screenshot shows the PubMed homepage at <https://www.ncbi.nlm.nih.gov/pubmed>. A search bar at the top contains the term "gastro". Below the search bar, a dropdown menu lists various search terms starting with "gastro", such as "gastric", "gastric cancer", "gastric bypass", etc. To the left of the search bar is a sidebar with links like "Using PubMed", "PubMed Quick Start Guide", and "Latest Literature". The main content area displays a list of recent articles, with the first few being "Blood (2)", "Cochrane Database Syst Rev (1)", "J Biol Chem (5)", and "J Clin Oncol (1)".

Using PubMed

[PubMed Quick Start Guide](#)

[Full Text Articles](#)

[PubMed FAQs](#)

[PubMed Tutorials](#)

[New and Noteworthy](#)

Latest Literature

New articles from highly accessed journals

Blood (2)

Cochrane Database Syst Rev (1)

J Biol Chem (5)

J Clin Oncol (1)

gastro

gastric

gastric cancer

gastric bypass

gastric ulcer

gastric emptying

gastric carcinoma

atrophic gastritis

early gastric cancer

cancer gastric

gastric acid

flat gastric

gastric sleeve

gastric varices

advanced gastric cancer

gastric banding

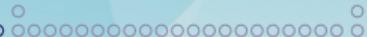
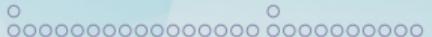
gastric bypass surgery

roux-en-y gastric bypass

gastric cancer review

chemotherapy gastric

helicobacter pylori gastric



PubMed as an IR engine

Screenshot of a web browser showing the PubMed search results for "adenocarcin*".

The search bar contains "adenocarcin*".

Filters applied:

- Article types: Case Reports
- Text availability: Free full text
- Species: Humans

Sort by: Most recent (selected)

Search results:

Items: 1 to 20 of 6790

1. A case report of a giant appendiceal mucocele and literature review.
Motsumi MJ, Motlaleslelo P, Ayane G, Sesay SO, Valdes JR.
Pan Afr Med J. 2017 Oct 4;28:106. doi: 10.11604/pamj.2017.28.106.13832. eCollection 2017. Review.

PMID: 29515724 [Free PMC Article](#)

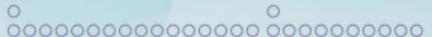
2. Concomitant endometrial and cervical adenocarcinoma: A case report and literature review.
Xu M, Zhou F, Huang L.
Medicine (Baltimore). 2018 Jan;97(1):e9596. doi: 10.1097/MD.00000000000009596. Review.

PMID: 29505548 [Free Article](#)

Results by year chart (Download CSV available).

Titles with your search terms:

- Single-incision laparoscopic surgery for loc: [Acta Gastroenterol Belg. 2018]
- The effectiveness of cycloox₁ [Biomed Pharmacother. 2018]
- Loss of NF-κB1 Causes Gastric



PubMed as an IR engine

Advanced search - PubMed + | X

https://www.ncbi.nlm.nih.gov/pubmed/advanced

NCBI Resources How To Sign in to NCBI

PubMed Home More Resources Help

YouTube Tutorial

PubMed Advanced Search Builder

drug induced liv[MeSH Terms]

Edit Clear

Builder

MeSH Terms

AND All Fields

Search or Add to history

History

There is no recent history

drug induced liv

chronic drug induced liver injury

disease, drug induced liver

diseases, drug induced liver

drug induced liver disease

drug induced liver diseases

drug induced liver injuries

drug induced liver injury

drug induced liver injury, chronic

injuries, drug induced liver

injury, drug induced liver

Show index list

Show index list

You are here: NCBI > Literature > PubMed

Support Center

GETTING STARTED

NCBI Education

NCBI Help Manual

NCBI Handbook

RESOURCES

Chemicals & Bioassays

Data & Software

DNA & RNA

POPULAR

PubMed

Bookshelf

PubMed Central

FEATURED

Genetic Testing Registry

PubMed Health

GenBank

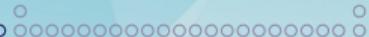
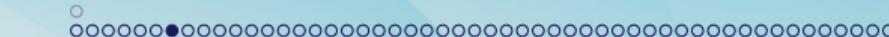
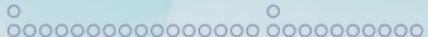
NCBI INFORMATION

About NCBI

Research at NCBI

NCBI News & Blog





PubMed as an IR engine

PubMed Advanced Search Builder



```
(((((drug induced liver injury[MeSH Terms]) AND "english"[Language]) AND "case reports"[Publication Type]) NOT surgery[Text Word]) OR administration research, nursing[MeSH Terms]) AND ("2000/01/01"[Date - Completion] : "3000"[Date - Completion])
```

[Edit](#)[Clear](#)

Builder

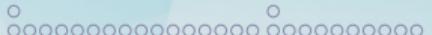
MeSH Terms	drug induced liver injury	Show index list	
AND	Language	"english"[Language]	Show index list
AND	Publication Type	"case reports"[Publication Type]	Show index list
NOT	Text Word	surgery	Show index list
OR	MeSH Terms	administration research, nursing	Show index list
AND	Date - Completion	2000/01/01 to present	Show index list
AND	All Fields		Show index list

[Search](#) or [Add to history](#)

History

There is no recent history





PubMed as an IR engine

(((drug induced liver injury[M... X | + https://www.ncbi.nlm.nih.gov/pubmed?term=((drug%20induced%20liver%20injury[MeSH%20Terms])%20AND%20(%222000/01/01%5BDate%20-%20Completion%5D%20%3A%20%223000%5BDate%20-%20Completion%5D))| Search NCBI Resources How To Sign in to NCBI Help

PubMed nursing[MeSH Terms]) AND ("2000/01/01"[Date - Completion] : "3000"[Date - Completion]) | Search

Create RSS Create alert Advanced

Article types Clinical Trial Review Customize ...

Text availability Abstract Free full text Full text

Publication dates 5 years 10 years Custom range...

Species Humans Other Animals

Clear all Show additional filters

Format: Summary Sort by: Most Recent Per page: 20 Send to + Filters: Manage Filters

Sort by: Best match Most recent

Items: 1 to 20 of 3158 << First < Prev Page 1 of 158 Next > Last >>

Search results

1. [Fimasartan-induced liver injury in a patient with no adverse reactions on other types of angiotensin II receptor blockers: A case report.](#)
Park DH, Yun GY, Eun HS, Joo JS, Kim JS, Kang SH, Moon HS, Lee ES, Lee BS, Kim KH, Kim SH.
Medicine (Baltimore). 2017 Nov;96(47):e8905. doi: 10.1097/MD.00000000000008905.
PMID: 29382024 Free PMC Article
[Similar articles](#)

2. [A rare case of methimazole-induced cholestatic jaundice in an elderly man of Asian ethnicity with hyperthyroidism: A case report.](#)
Ji H, Yue F, Song J, Zhou X.
Medicine (Baltimore). 2017 Dec;96(49):e9093. doi: 10.1097/MD.00000000000009093.
PMID: 29245333 Free PMC Article
[Similar articles](#)

3. [Copper-associated hepatitis in a patient with chronic myeloid leukemia following hematopoietic stem cell transplantation: A case report.](#)
Lee CF, Chen CH, Wen YC, Chang TY, Lai MW, Jaing TH.
Medicine (Baltimore). 2017 Dec;96(49):e9041. doi: 10.1097/MD.00000000000009041.
PMID: 29245301 Free PMC Article
[Similar articles](#)

Results by year

Download CSV

Titles with your search terms

Nursing Administration Research Priorities: Findings From a Delphi Study [J Nurs Adm. 2016]

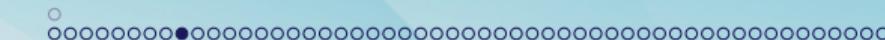
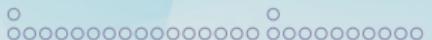
Building capacity for the conduct of nursing research at a Veterans Administration [J Nurs Adm. 2015]

Nursing administration research: an evolving science. [J Nurs Adm. 2014]

See more...

Find related data

Database: Select

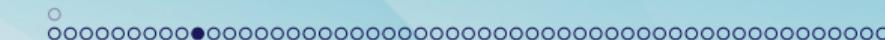
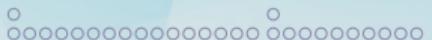


PubMed as an IR engine

The screenshot shows a web browser window for 'Details - PubMed - NCBI'. The address bar shows the URL: <https://www.ncbi.nlm.nih.gov/pubmed/details?term=lung+cancer>. The search term 'lung cancer' is entered in the search bar. The results page displays the following sections:

- Search Details:**
 - Query Translation:**

```
"lung neoplasms"[MeSH Terms] OR ("lung"[All Fields] AND "neoplasms"[All Fields]) OR "lung neoplasms"[All Fields] OR ("lung"[All Fields] AND "cancer"[All Fields]) OR "lung cancer"[All Fields]
```
 - Result:** 302129
- Translations:**
 - lung "lung neoplasms"[MeSH Terms] OR ("lung"[All Fields] AND "neoplasms"[All Fields]) OR "lung neoplasms"[All Fields] OR ("lung"[All Fields] AND "cancer"[All Fields]) OR "lung cancer"[All Fields]
- Database:** PubMed
- User query:** lung cancer



PubMed as an IR engine

Details - PubMed - NCBI

https://www.ncbi.nlm.nih.gov/pubmed/details/?

NCBI Resources How To Sign in to NCBI Help

PubMed pten mutation Advanced

Search Details

Query Translation:

```
pten[All Fields] AND ("mutation"[MeSH Terms] OR "mutation"[All Fields])
```

Search URL

Result:

4194

Translations:

```
mutation ("mutation"[MeSH Terms] OR "mutation"[All Fields])
```

Database:

PubMed

User query:

pten mutation

Tobramycin (D014031)

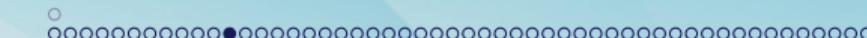
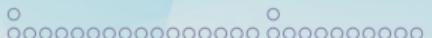
Gentamicins (D005839)

We observed patients treated with gentamicin sulfate or tobramycin sulfate for the development of aminoglycoside-related renal failure. Gentamicin sulfate decreased renal function more frequently than tobramycin sulfate.

Aminoglycosides (D000617)

Renal Insufficiency (D051437)

MeSH Indexing



[J Alzheimers Dis.](#) 2012;32(2)

Disease

A

Gene

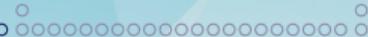
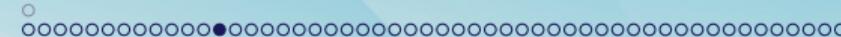
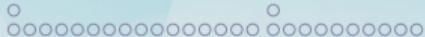
Variant

Highly pathogenic Alzheimer's disease presenilin 1 P117R mutation causes a specific increase in p53 and p21 protein levels and cell cycle dysregulation in human lymphocytes.

[Bialopiotrowicz E¹](#), [Szybinska A](#), [Kuzniewska B](#), [Buizza L](#), [Uberti D](#), [Kuznicki J](#), [Wojda U](#).

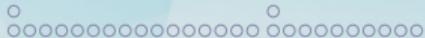


Rank	Gene or Protein ID	Gene SYM	WTAA	MTAA	POS	Disease	PMIDs
1	Q13131	PRKAA1	Q	R	16	Breast cancer	16959974
2	P31749	AKT1	E	K	17	Breast cancer	17611497 18954143 19713527 21793738
3	P10275	AR	H	Y	874	Prostate cancer	17591767



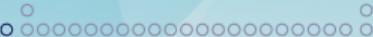
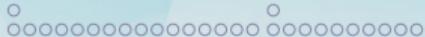
NLP steps

- **sentence segmentation** splitting a multi-sentence input into its individual sentences
in biomedical texts, sentences can even begin with lower-case letters, e.g.,
lush mutants are also defective for pheromone-evoked behavior (PMID 15664171)



NLP steps

- **Tokenization** is the task of splitting an input into individual words (as well as other units, such as punctuation marks).
 - apparently trivial, but in fact difficult even in general text
 - plus a number of domain-specific challenges in biomedical text Example: hyphen.
 - *-fever* absence of a fever
 - *Cl-* negative charge on a chlorine ion
 - *BRCA1-* non-functional gene
 - *BRCA1-IRIS* joins two synonymous forms of the gene name
 - *BRCA1-mediated* syntactic, not semantic
 - *BRCA1-CtIP* the BRCA1 protein forms a complex with the CtIP substrate of ATM protein kinase



NLP steps

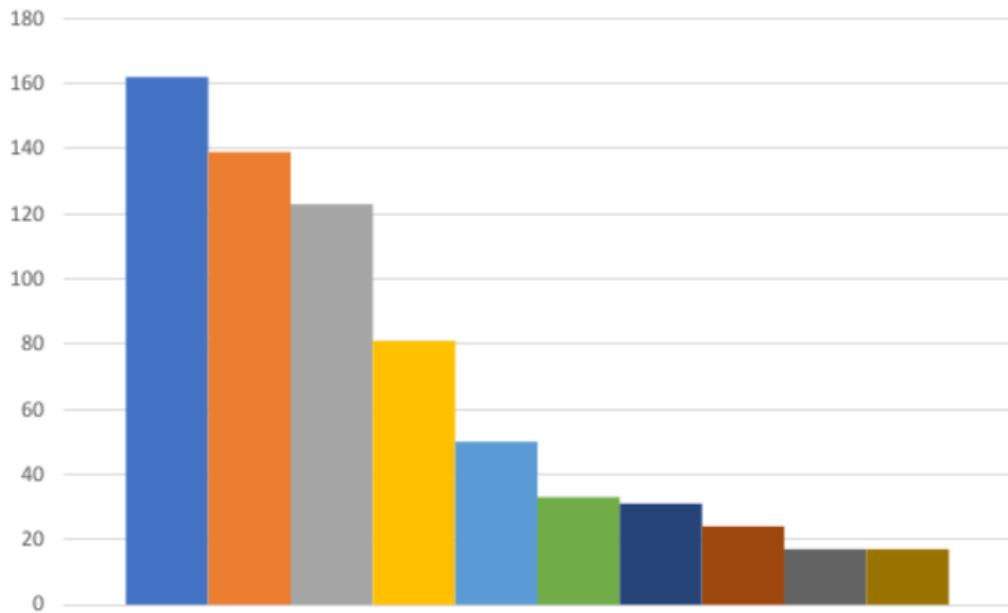
- **Abbreviations** E.g. does "PDA" mean *patent ductus arteriosus* or *posterior descending artery* or *phorbol 12,13 diacetate*?

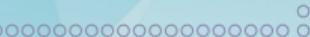
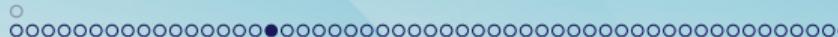
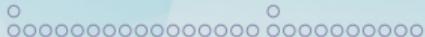
Almost 22% of abbreviations in one sample of biomedical text having more than one possible expansion and an average of 4.61 possible definitions for abbreviations six or fewer characters long [[Chang et al., 2002](#)]

Possible approach: Schwartz and Hearst's algorithm (Schwartz and Hearst 2003); easily implementable; returns a single output with a high likelihood of being correct, and of requiring no connections with external applications



- superficial capillary plexus
- superior cerebellar peduncle
- slow cortical potentials
- survivorship care plan
- single cell protein
- supercooling point
- selective cerebral perfusion
- sperm coating protein
- sulfachloropyridazine
- survival continence and potency





NLP steps

- **Named Entity Recognition**

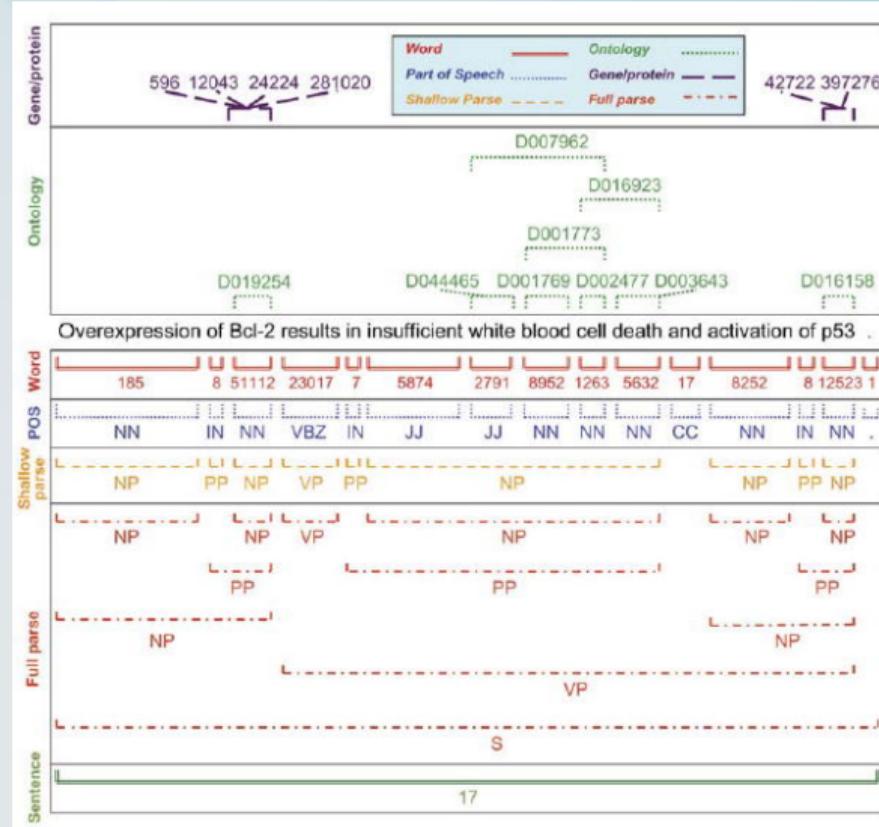
- Entity recognition is a required component for almost any biomedical text mining application.
- Examples
 - ABNER system (Settles 2005). Its advantages include robustness (it is distributed as a .jar file and works well on multiple platforms) and a variety of input and output format options. In addition to gene/protein names, it also locates cell types and cell lines.
 - LingPipe's gene-mention-finding module
 - MetaMap tool (Aronson 2001) recognition of a wider range of biomedically relevant categories—diseases, drugs, chemicals, treatment modalities, and the like, available as a Java API or via a server. Recognizes variant forms of domain-relevant concepts
 - Banner <http://banner.sourceforge.net/>

O
ooooooooooooooooooo

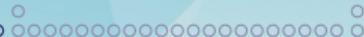
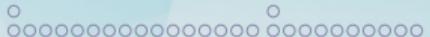
O
oooooooooooooooooooo

●
oooooooooooooooooooo

O
oooooooooooooooooooo

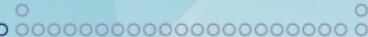
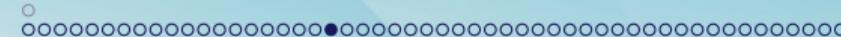
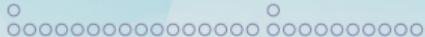


BioNLP analysis levels
[Hunter and Cohen, 2006]



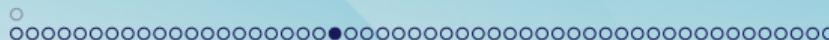
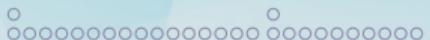
TM approaches

Approach	Description	Examples
Manual Extraction	Literature curation, high precision Domain expertise, small scale	GO annotation, COSMIC, LTKB, Drugbank, SwissProt,..
Dictionaries, gazetteers	Match entity/term lists. Document indexing. Term co-occurrence	Location: membrane, mitochondria, cytosol, nucleus, ..
Pattern based	Recurrent word patterns & language expressions.	<CHEMICAL>V:induced<TERM:HEPATOBILIARY> <PROTEIN>and<PROTEIN>interact
Knowledge based (rule)	Knowledge about language structure, order, POS, triggers	RLIMS-P [<AGENT> phosphorylation]NP of <THEME>
Machine learning	Labeled training, building classifiers, word features	ABNER, MedlineRanker, CheNER
Hybrid systems	Incorporate statistical model, dictionaries and patterns/rules	tmChem, Chemspot

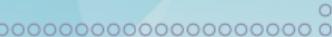
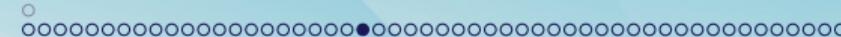
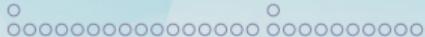


Problems

- **Synonyms and Term variability: one concept, multiple names** For example, BRCA1 could be referred to by any of its alternate symbols *IRIS*, *PSCP*, *BRCAI*, *BRCC1*, or *RNF53* (or by any of their many spelling variants, which include *BRCA1*, *BRCA-1*, and *BRCA 1*) or by any of the variants of its full name, viz. *breast cancer 1, early onset* (its official name per Entrez Gene and the Human Gene Nomenclature Committee), as *breast cancer susceptibility gene 1*, or as the latter's variant *breast cancer susceptibility gene-1*. Similarly, *breast cancer* could be referred to alternatively as *carcinoma of the breast*, or *mammary neoplasm*.

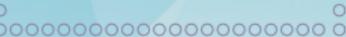
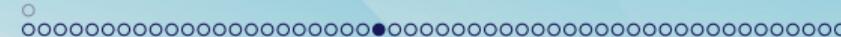
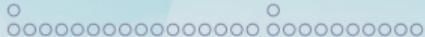


- pneumonia of unknown aetiology
- 2019-nCov infection
- novel coronavirus pneumonia
- SARS-CoV-2 infection
- COVID-19
- coronavirus disease 2019
- 2019 novel coronavirus infection disease
- nCOVID-19
- severe acute respiratory syndrome coronavirus 2 infection
- Coronavirus disease of 2019
- Wuhan coronavirus pneumonia
- CoV 19 infection
- COIVD-19 disease
- COVID-19 ARDS
- nCov-19 infection
- SARS-CoV-2 infectious disease
- CV-19
- coronavirus 2 syndrome
- SARS-CoV-2 associated ARDS
- C19



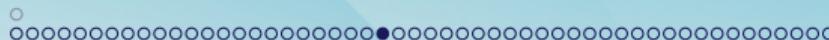
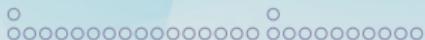
Problems

- **Ambiguity: one name, multiple concepts/meanings** A primary problem that either type of system must deal with is the issue of ambiguity: the existence of multiple relationships between language and meanings or categories. Ambiguity exists at every level of linguistic structure, from the part of speech of words to subtle issues in pragmatics.
 - with general language
 - across entity types
 - metonymy
 - within the same entity type
 - across species
 - within a species
 - grammatical ambiguity



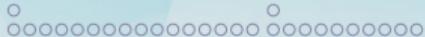
Ambiguity

- **across species** fat also turns out to be the name or symbol of a number of different genes. Humans, mice, rats, Drosophila, zebrafish, chickens, M. mulatta, and two Lactobacilli have at least one gene whose name, official symbol, or alias is fat.
- **within a species** in humans, fat is the official symbol of Entrez Gene entry 2195 and an alternate symbol for Entrez Gene entry 948. The distinction is not trivial. The former is a cadherin, and is associated with tumor suppression and with bipolar disorder, while the latter is a thrombospondin receptor associated with atherosclerosis, platelet glycoprotein deficiency, hyperlipidemia, and insulin resistance, to name just a few phenotypes. These ambiguities are not trivial: if your analysis is wrong, you miss or erroneously extract information on relations between molecular biology and human disease.



- "... suggests a role for BRCA1 in HIF-1alpha regulation."
 - Human = NCBI Gene 672
 - M. musculus = NCBI Gene 12189
 - R. norvegicus = NCBI Gene 497672
- "... to lower blood cholesterol levels in comparison to sugar cane policosanol (SCP) in rabbits."

Source: PMID 18030615



Ambiguity

- **grammatical ambiguity** *fat*: noun or verb? PubMed returns almost 112 K hits for that single-word query (and more than 13 K even if we try to restrict the query to genomics by including the disjunction (gene OR genetic OR genetics)).
- **using a common term** many Drosophila symbols and names are the same as common English words: *a*, *to*, and *And* are all symbols of Drosophila genes (Entrez Gene IDs 43852, 43036, and 44913, respectively) [[Hunter and Cohen, 2006](#)]
- **metonymy** “referring to something by an entity that is related to it”: it is often not clear whether a string like p53 refers to the gene of that name, to the protein that it codes for, or to its mRNA



Met28 binds to DNA

...binding of Met28 to DNA...

...Met28 and DNA bind...

...binding between Met28 and DNA...

...Met28 is sufficient to bind DNA...

...DNA bound by Met28...



Generic Environments for TM pipelines

- GATE (General Architecture for Text Engineering, Cunningham 2002) system. It has been available since the mid 1990s, is used by many groups, and has been heavily field-tested over the years.
- UIMA (Unstructured Information Management Architecture, Ferrucci and Lally 2004, Mack et al. 2004), IBM, java
- LingPipe, produced by Alias-i, offers a mid-point of sorts between un-integrated pipelines and comprehensive architectures.

GATE Developer 6.1-snapshot build 1911

Annotation Sets **Annotations List** **Annotations Stack** **Co-reference Editor** **OAT** **RAT-C** **RAT-I** **Text**

Annotations: water_goes_out... gaze 12B.xml_003D

The HUGO Gene Nomenclature Committee (HGNC) provides a list of genes in the WNT gene family.

Genetics Home Reference summarizes the normal function and health implications of these members of the WNT gene family: WNT3 and WNT4.

What conditions are related to genes in the WNT gene family?

Genetics Home Reference includes these conditions related to genes in the WNT gene family:

- WNT4 Mütterian aplasia or dysplasia
- WNT4 Müllerian aplasia or dysplasia dysfunction
- Where can I find additional information about the WNT gene family?

You may find the following resources about the WNT gene family helpful:

- The Wnt Homepage, Stanford University
- Developmental Biology (Sefti et al., 2005). The Wnt signal transduction pathway (figure)
- Molecular Biology of the Cell (Lions et al., 2002). Wnt Proteins Bind to Frizzled Receptors and Inhibit the Degradation of Beta-Catenin
- List of human WNT genes, Stanford University
- Where can I find general information about genes and gene families?

The Handbook provides basic information about genetics in clear language.

What is DNA?
 What is a gene?
 What are gene families?
 What glossary definitions help with understanding the WNT gene family?

birth defect ; defect ; cell ; degenerative ; differentiation ; embryonic ; gene ; ligand ; molecule ; mutation ; nerve cell ; nucleus ; proliferation ; protein ; receptor ; stage ; stem cells ; tissue

You may find definitions for these and many other terms in the Genetics Home Reference Glossary.

See also Understanding Medical Terminology.

References (8 links)

The resources on this site should not be used as a substitute for professional medical care or advice. Users seeking information about a personal genetic disease, syndrome, or condition should consult with a qualified healthcare professional. See How can I find a genetics professional in my area? in the Handbook.

Reviewed: February 2009
 Published: February 27, 2011
 URL: <http://ghr.nlm.nih.gov>, National Center for Biomedical Communications
 U.S. National Library of Medicine, National Institutes of Health
 Department of Health & Human Services, USA.gov
 Freedom of Information Act Copyright Privacy Accessibility
 Indicates a page outside Genetics Home Reference.
 Links to web sites outside the Federal government do not constitute an endorsement.
 See Selection Criteria for Web Links.
 This site complies with the HONcode standard for trustworthy health information: verify here.

Document Editor Initialisation Parameters

Online help for the selected component

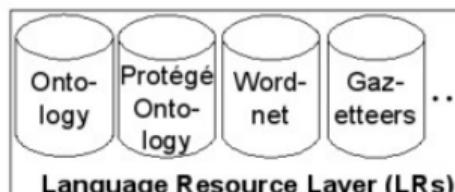
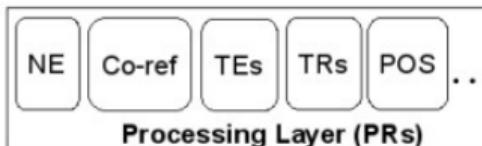
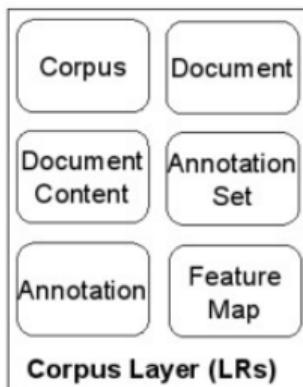
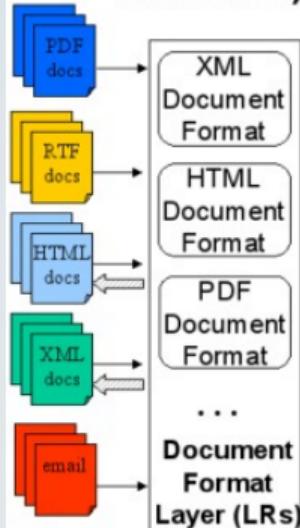
Annotations: water_goes_out... gaze 12B.xml_003D

MISCCELLANEOUS

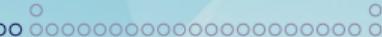
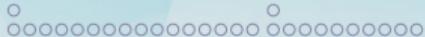
- Molecular_Function
- Neoplastic_Process
- Nucleic_Acid_Nucleoside_or_Nucleotide
- Occupation_or_Disipline
- Occupational_Activity
- Organ_or_Tissue_Function
- Organic_Chemical
- Organism
- Organism_Attribute
- Organism_Function
- Organization
- Page
- Pharmacologic_Substance
- Phenomenon_or_Process
- Physiologic_Function
- Plant
- Population_Group
- Professional_or_Occupational_Group
- Qualitative_Concept
- Quantitative_Concept
- Receptor
- Regulation_or_Law
- Research_Activity
- Section
- Self_Help_or_Relief_Organization
- Sentence
- Sign_or_Symptom
- Social_Behavior
- SpaceToken
- Spatial_Concept
- Split
- Substance
- Temporal_Concept
- Tissue
- Token
- Virus
- Original_markups
- Ifthen
- Anatomy
- Problem

GATE

APIs (GATE Embedded)

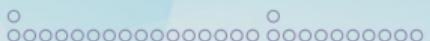


- NOTES**
- everything is a replaceable bean
 - all communication via fixed APIs
 - low coupling, high modularity, high extensibility



Ontologies

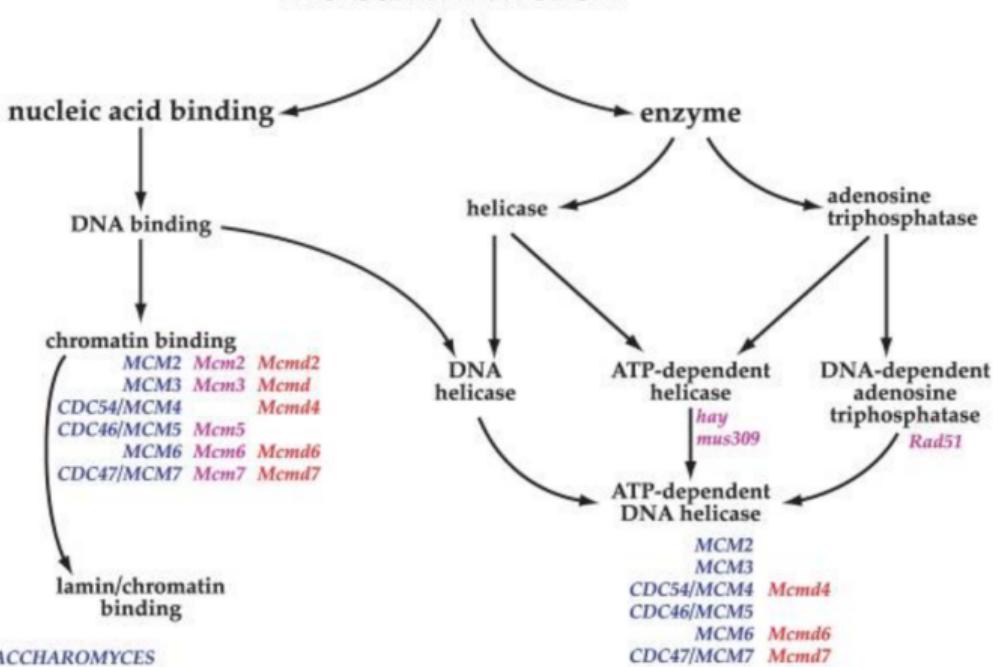
- are the most effective representational mechanism
- verbs and relations might benefit from richer, frame-like or event-related structure
- Gene Ontology [[Ashburner et al., 2000](#)] The ontologies are built from a structured, controlled vocabulary.

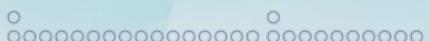


GO

b

molecular function

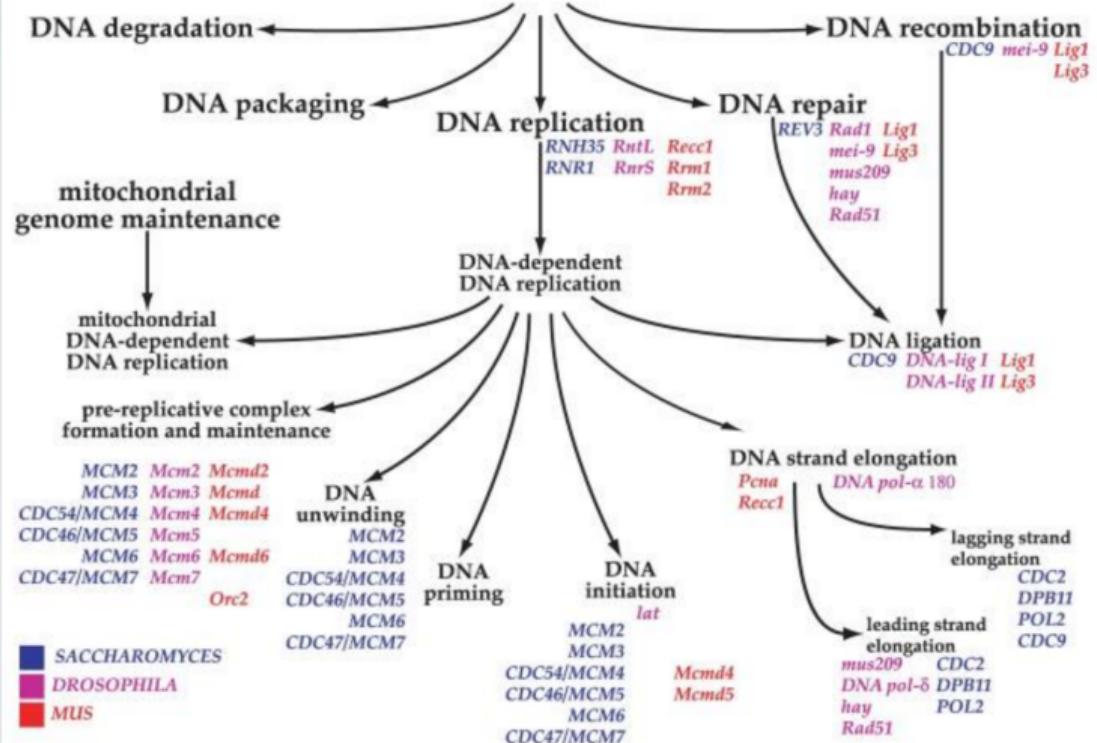


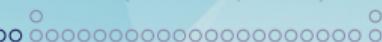
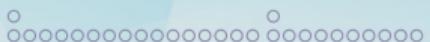


GO

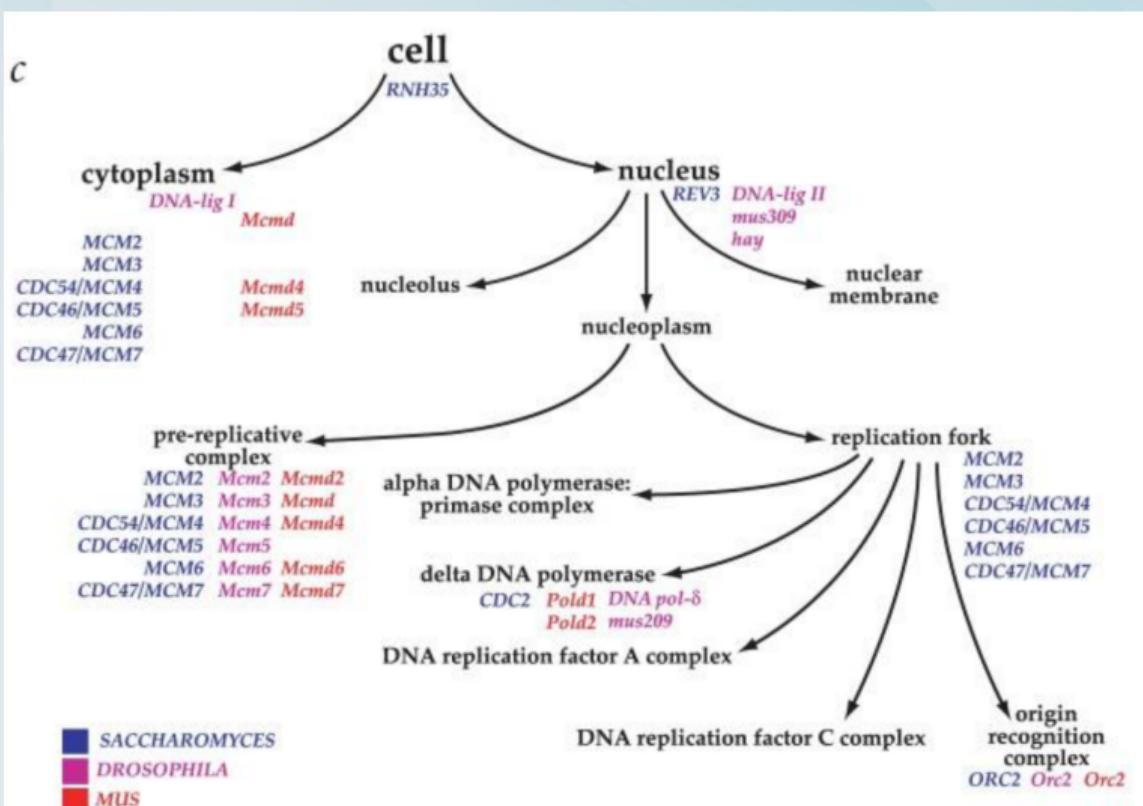
a

DNA metabolism



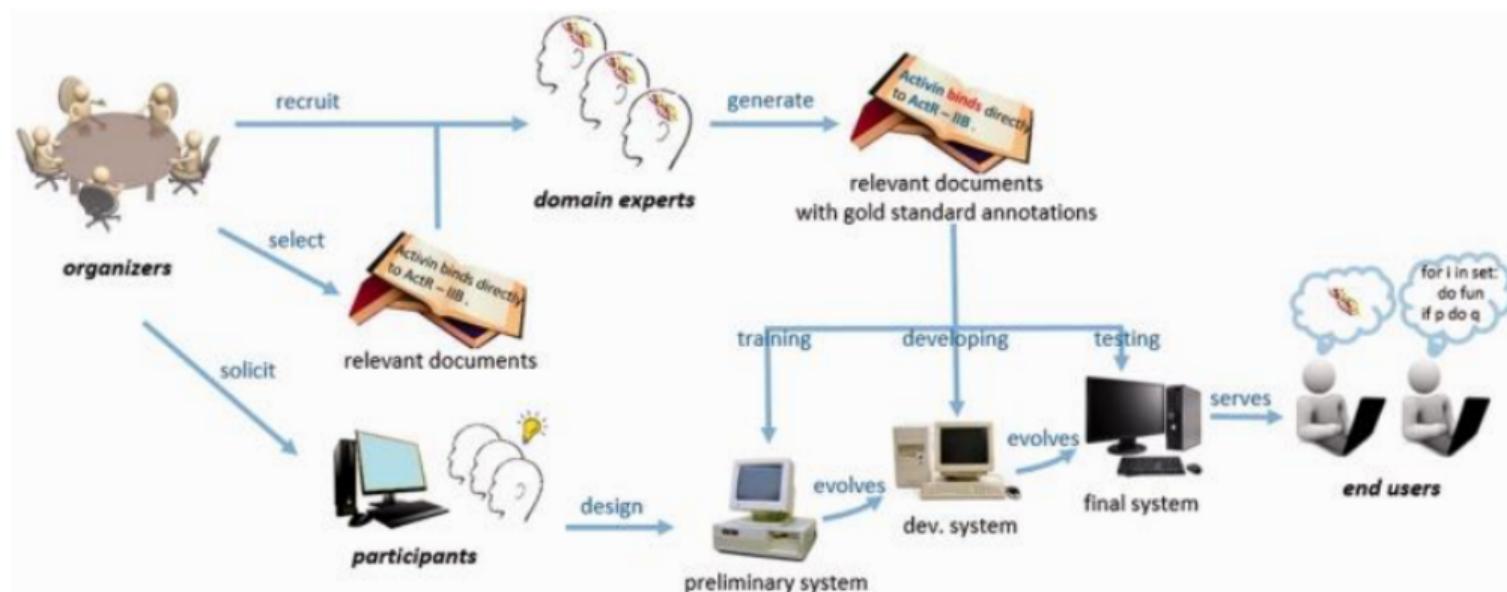


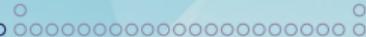
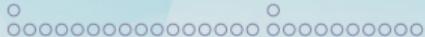
GO



Competitive Evaluations

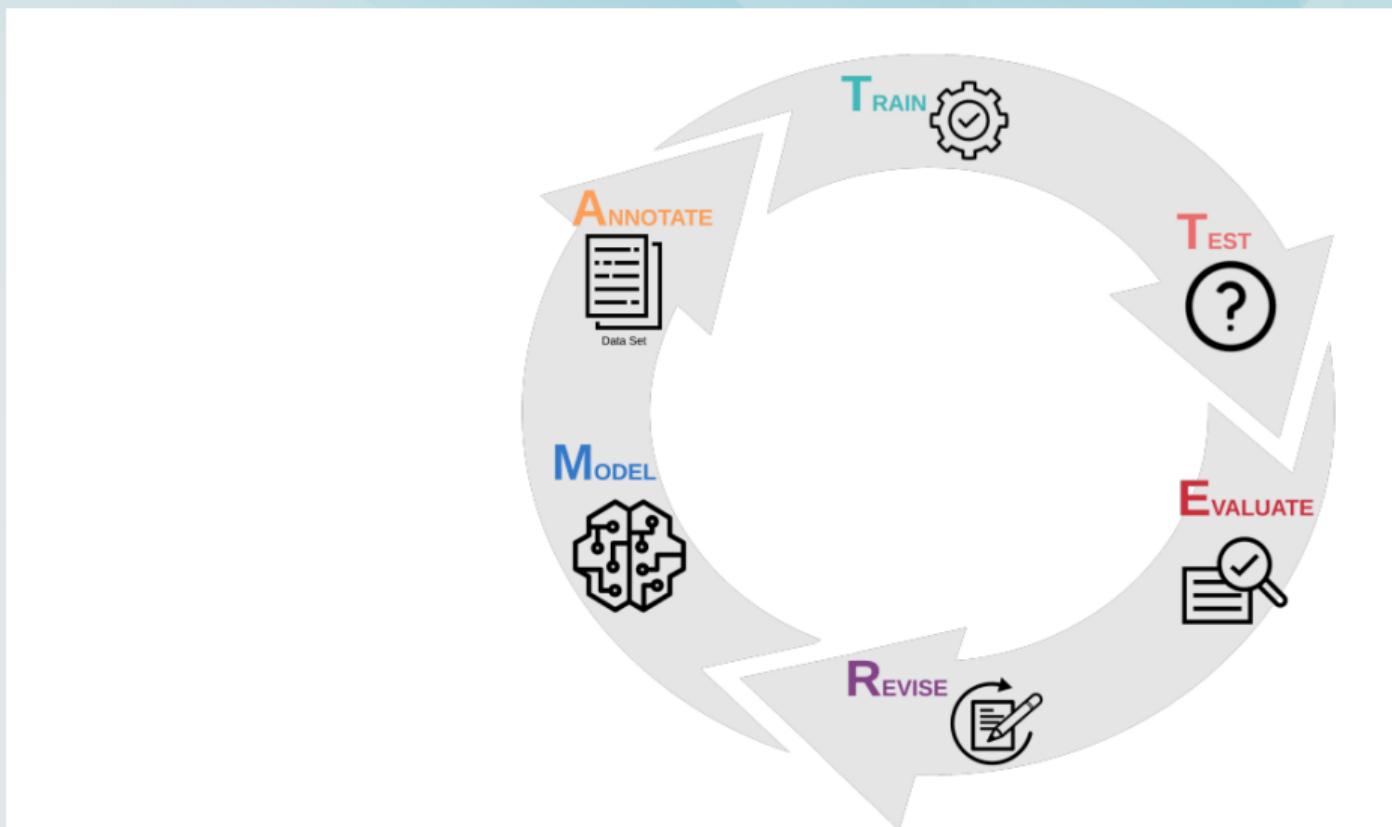
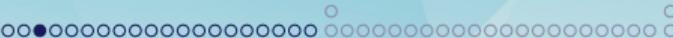
- also called "Shared Tasks", "Competitive scientific challenges"

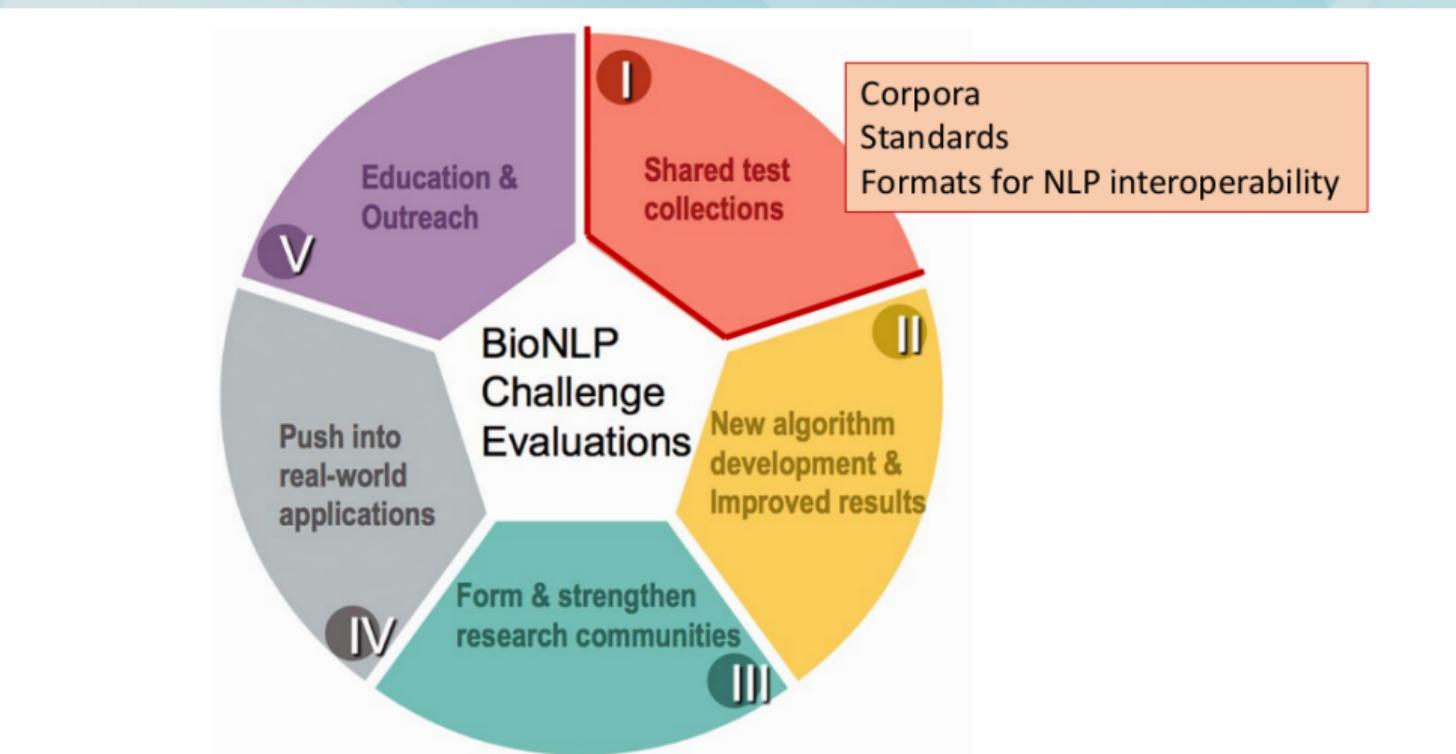
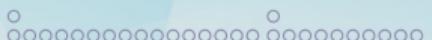


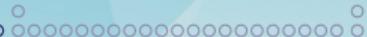
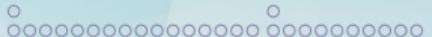


Competitive Evaluations

- TREC Genomics Track (2003-2007) <http://trec.nist.gov/data/genomics.html>
Supported by the National Institute of Standards and Technology, TREC provides a forum for evaluation of information retrieval systems.
- BioNLP competition
- BioCreative
- IE Task @ KDD Cup earliest biomedical text data mining competition, an information extraction task sponsored by the Knowledge Discovery and Data mining (KDD) Cup in which participants built systems to aid in the FlyBase curation process (Yeh et al., 2003)
[Yeh et al., 2003]

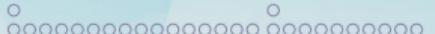




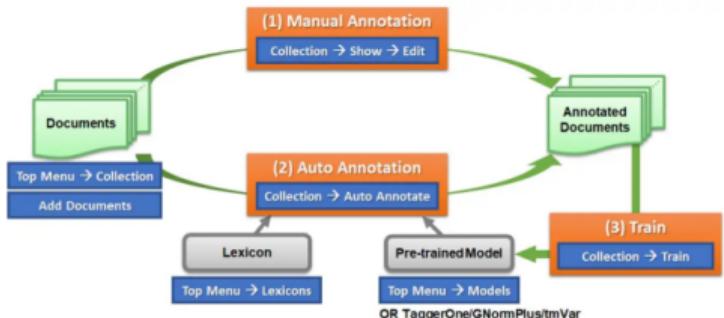


<http://www.biocreative.org>

- High quality corpora
- Covers abstracts, full text and patents
- Annotation guidelines aligned with domain experts
- Inter-annotator agreement



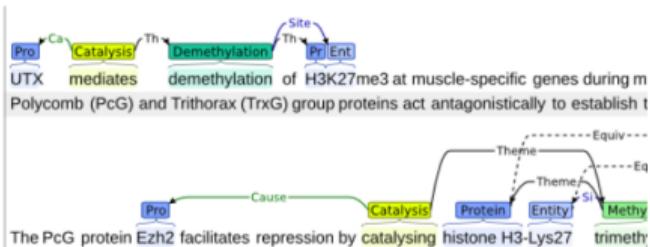
ezTag <https://eztag.bioqrator.org/>



TeamTat <https://www.teamtat.org/>

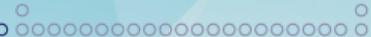
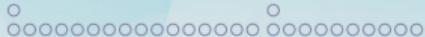
The screenshot shows a search result for 'Acetylation of the Cell-Fate Factor Dachshund Determines p53 Binding and Signaling Modules in Breast Cancer'. The results table includes columns for ID, Type, and Nodes. A detailed abstract is provided, mentioning authors like Chen, Ke, Wu, Kongming, Gammie, Michael, Ertel, Adam, Wang, Jing, Zhang, Wei, Zhou, Ke, D'Silva, Gabriele, Li, Zhiping, Rui, Holstein, Qiong, Andrew, Alvarado, and others, along with journal information (Oncogene, 2013, Vol. 32, issue 6, 923–935). A note indicates that the gene was cloned as an inhibitor of the receptor tyrosine kinase growth factor (EGFR) ellipse.

brat <https://brat.nlplab.org/>



TextAE <http://textae.pubannotation.org/>

The screenshot shows a search result for 'EGFR production and NF-kappaB activation, which could both be blocked by antioxidants or EGFR inhibitors'. The results table includes columns for ID, Type, and Nodes. A detailed abstract is provided, mentioning authors like Chen, Ke, Wu, Kongming, Gammie, Michael, Ertel, Adam, Wang, Jing, Zhang, Wei, Zhou, Ke, D'Silva, Gabriele, Li, Zhiping, Rui, Holstein, Qiong, Andrew, Alvarado, and others, along with journal information (Oncogene, 2013, Vol. 32, issue 6, 923–935). A note indicates that the gene was cloned as an inhibitor of the receptor tyrosine kinase growth factor (EGFR) ellipse.



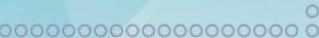
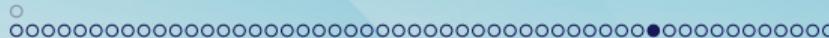
Literature-Based Discovery

Literature-based discovery is the attempt to automatically induce novel hypotheses by processing existing publications. [Swanson, 1988]

Early example: Arrowsmith system by Don R. Swanson (1924-2012).

"Imagine that the pieces of a puzzle are independently designed and created, and that, when retrieved and assembled, they then reveal a pattern undesigned, unintended, and never before seen, yet a pattern that commands interest and invites interpretation. So it is, I claim, that independently created pieces of knowledge can harbor an unseen, unknown, and unintended pattern. And so it is that the world of recorded knowledge can yield genuinely new discoveries"

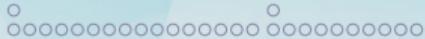
https://en.wikipedia.org/wiki/Arrowsmith_System



Chilibot

<http://www.chilibot.net> [Chen and Sharp, 2004]

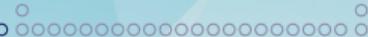
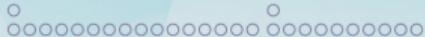
Chilibot is a freely available, web-based application that takes sets of gene names, and (optionally) additional keywords as input, and finds information about the relationships among them. In an initial information retrieval step, it constructs queries to send to the PubMed search engine. It uses basic language processing techniques to identify sentences that describe stimulatory, inhibitory, and other relationships between pairs of genes. For pairs of genes for which large sets of sentences are found, it then uses techniques from the field of automatic summarization (the area of natural language processing concerned with constructing shortened versions of input texts) to select the best sentences to display to the user. A graph displays the entire set of interactions amongst all of the genes that were input. By clicking on links between pairs of genes, the user can see the full set of sentences that describe relations between those genes, and clicking on the sentences themselves displays the original PubMed abstract.



Textpresso

<http://www.textpresso.org> [Miller et al., 2004]

The Textpresso system is an example of an information extraction project that strives to create an up-to-date summary of current knowledge related to a specific model organism. Users can search for information involving any of the 33 semantic classes of things and relationships between them that the system extracts information about (e.g., genes, clones, pathways, or mutations). The original system, devoted to *C. elegans*, has scanned thousands of full-text articles, and the resulting database is accessed more than 1500 times daily by nematode biologists worldwide. Work is in progress to extend this method to a variety of other organisms, including *N. crassa*, *D. melanogaster*, and *S. cerevisiae*, and preliminary systems are now available on several model organism database sites.



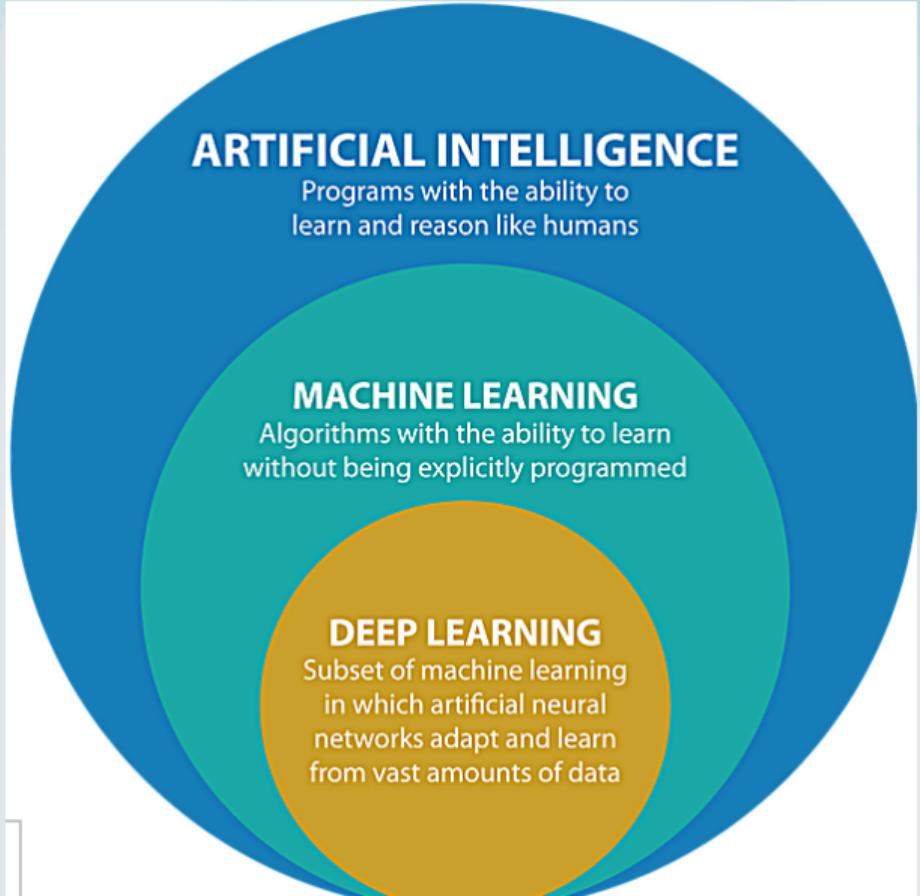
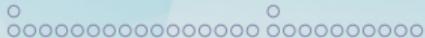
iHOP

<http://www.ihop-net.org/>

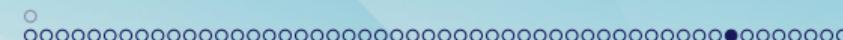
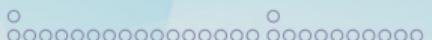
Information Hyperlinked over Proteins (or iHOP) is an online text-mining service that provides a gene-guided network to access PubMed abstracts. The service was established by Robert Hoffmann and Alfonso Valencia in 2004.

The concept underlying iHOP is that by using genes and proteins as hyperlinks between sentences and abstracts, the information in PubMed can be converted into one navigable resource. Navigating across interrelated sentences within this network rather than the use of conventional keyword searches allows for stepwise and controlled acquisition of information. Moreover, this literature network can be superimposed upon experimental interaction data to facilitate the simultaneous analysis of novel and existing knowledge

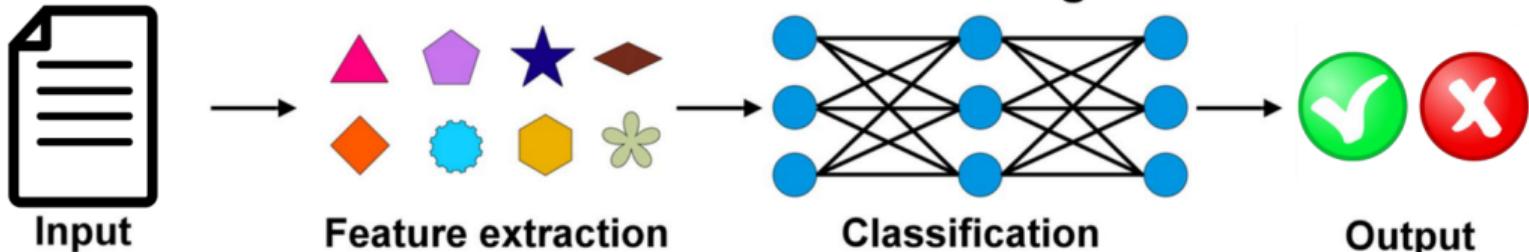
https://en.wikipedia.org/wiki/Information_Hyperlinked_over_Proteins



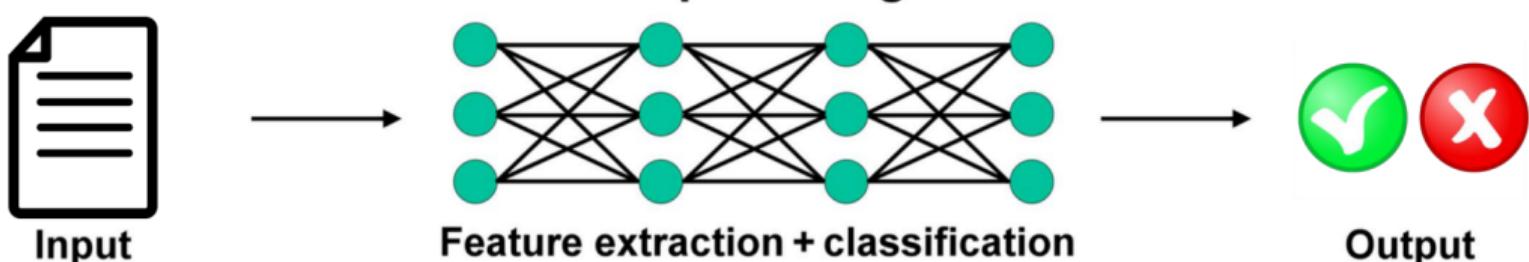
What is deep learning?



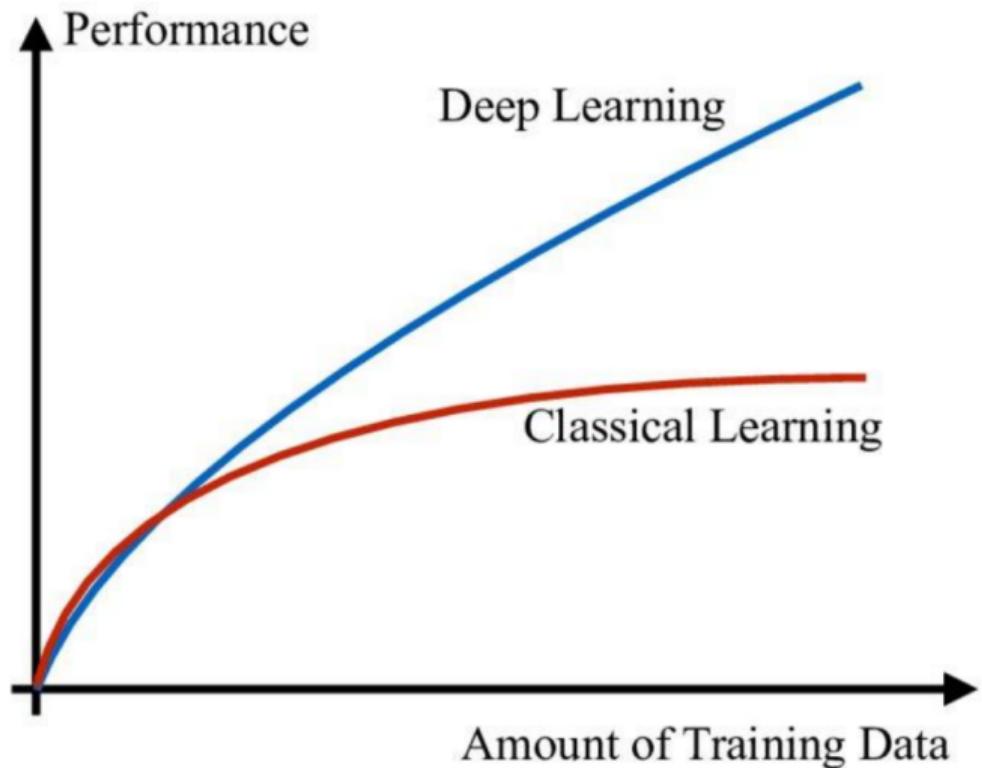
Traditional machine learning



Deep learning

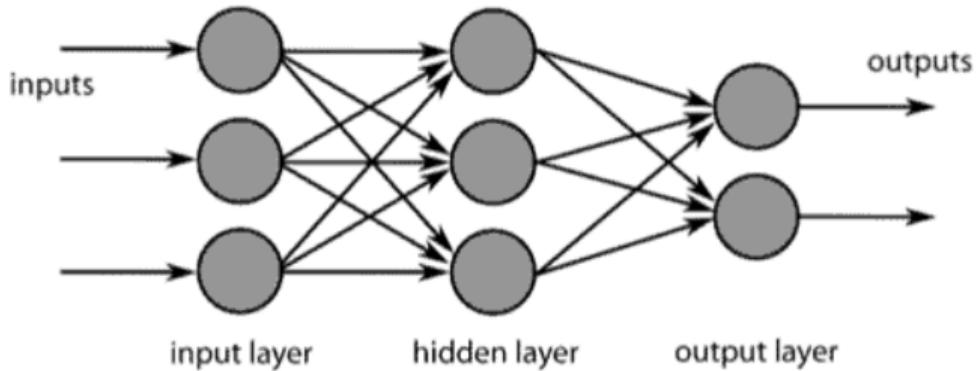


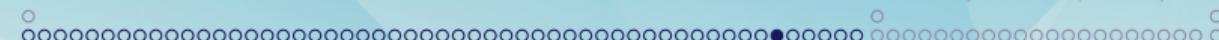
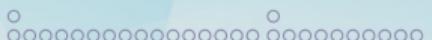
Renn, Alex, et al. "Advances in the prediction of mouse liver microsomal studies: From machine learning to deep learning." *Wiley Interdisciplinary Reviews: Computational Molecular Science* (2020): e1479.



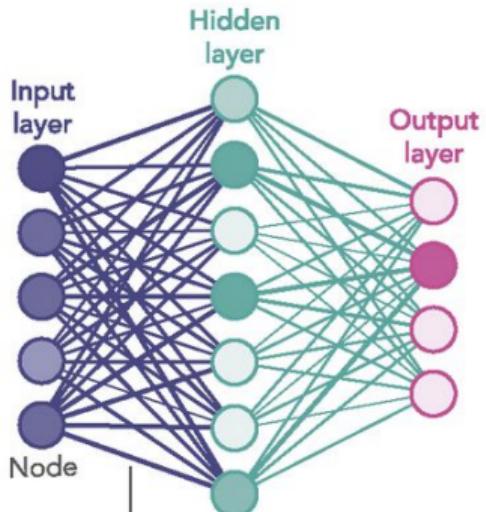
Zappone, Alessio, Marco Di Renzo, and Mérouane Debbah. "Wireless networks design in the era of deep learning: Model-based, AI-based, or both?." *IEEE Transactions on Communications* 67.10 (2019): 7331-7376.

Was ist Deep Learning?

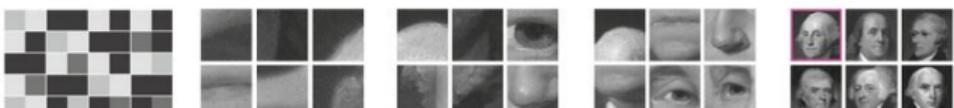
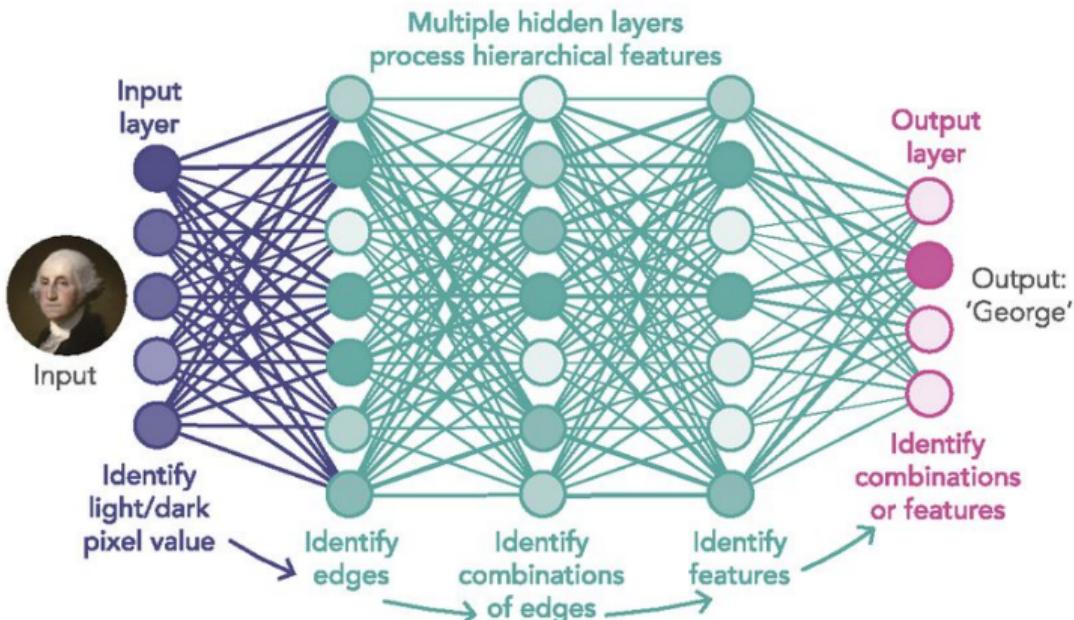


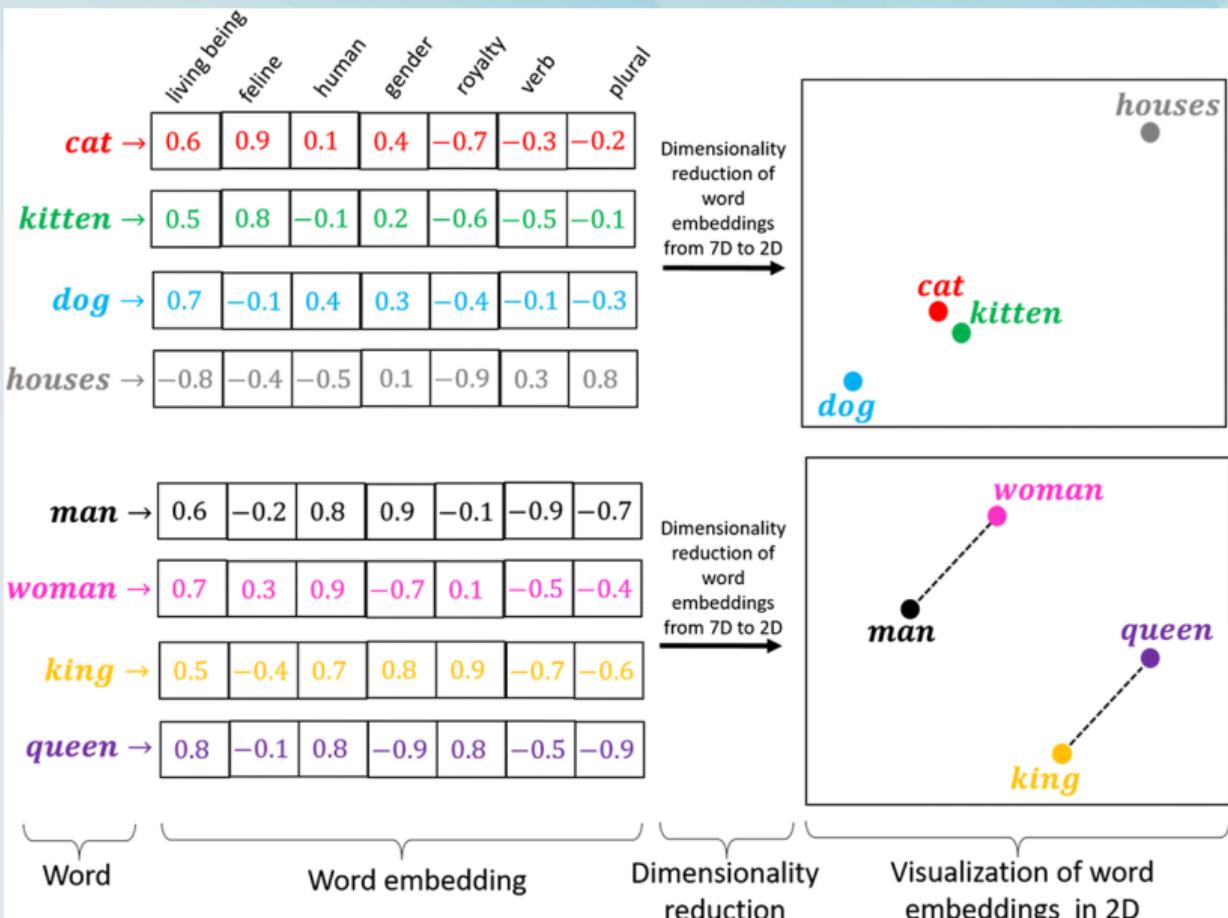
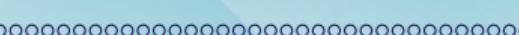
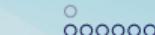
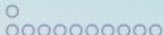


1980S-ERA NEURAL NETWORK



DEEP LEARNING NEURAL NETWORK



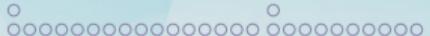




<i>word2vec</i>
<ul style="list-style-type: none">• Focus on local context• Implicitly learn embeddings based on local co-occurrences• Has two versions: cbow and skip-gram

<i>glove</i>
<ul style="list-style-type: none">• Focus on global context• Explicitly learn embeddings based on global co-occurrences

<i>fastText</i>
<ul style="list-style-type: none">• Extends from word2vec by using word subgrams• Predict an out-of-vocabulary word based on subgram statistics



ELMo
Oct 2017
Training:
800M words
42 GPU days

GPT
June 2018
Training
800M words
240 GPU days

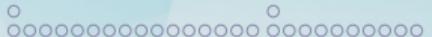
BERT
Oct 2018
Training
3.3B words
256 TPU days
~320–560
GPU days

GPT-2
Feb 2019
Training
40B words
~2048 TPU v3
days according to
[a reddit thread](#)

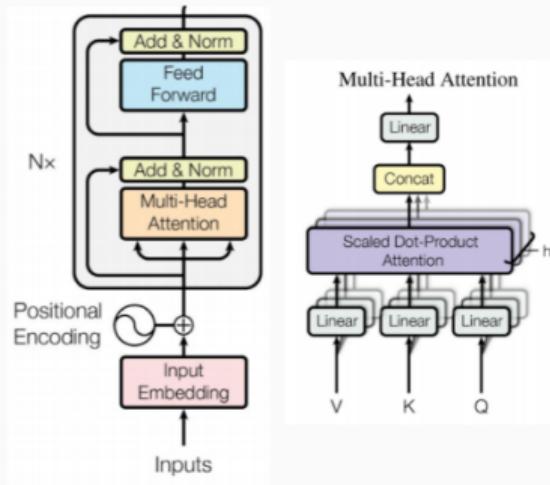
XL-Net,
ERNIE,
Grover
RoBERTa, T5
July 2019—



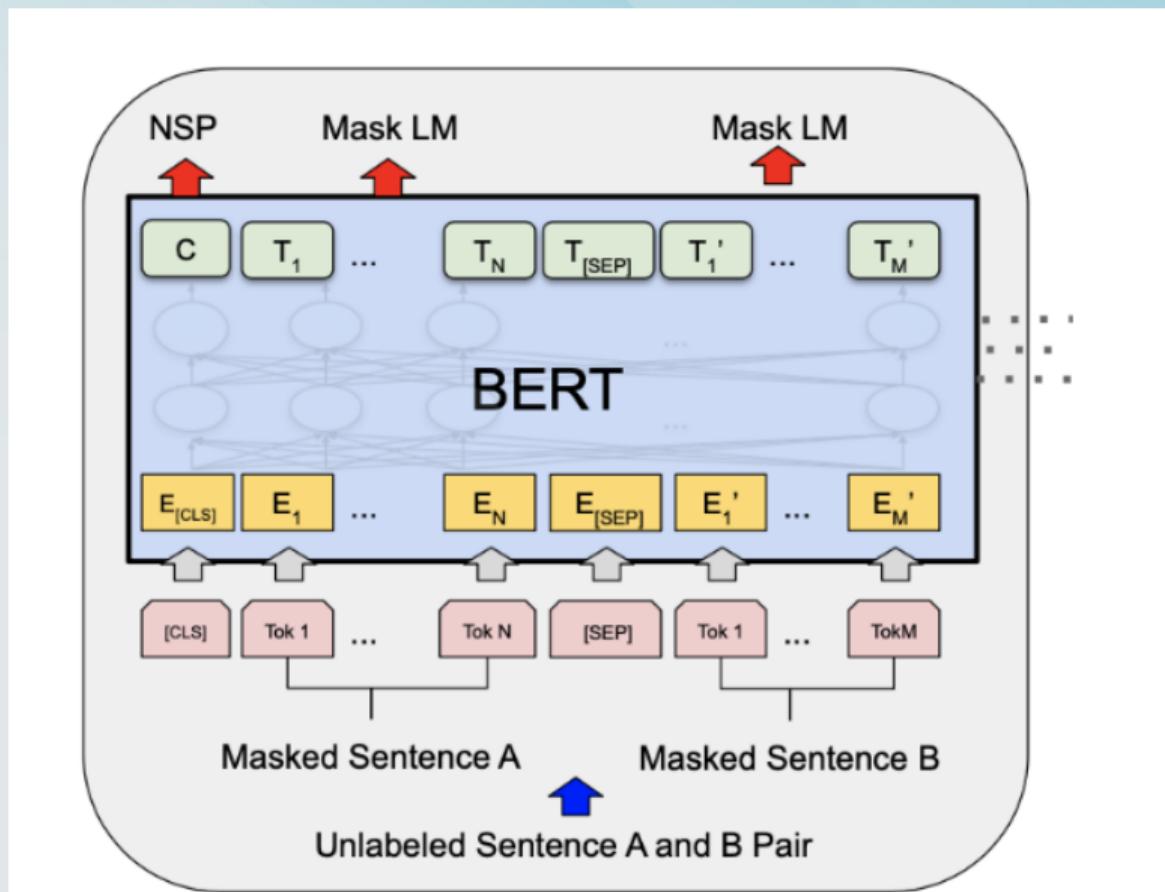
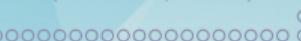
All of these models are Transformer models



- Multi-headed self attention
 - Models context
- Feed-forward layers
 - Computes non-linear hierarchical features
- Layer norm and residuals
 - Makes training deep networks healthy
- Positional embeddings
 - Allows model to learn relative positioning



<http://jalammar.github.io/illustrated-transformer/>



Introduction

The Biomedical Literature Methods



Our Tools ("OntoGene / IDSIA")

References



Topic

Introduction

The Biomedical Literature

Methods

Our Tools ("OntoGene / IDSIA")

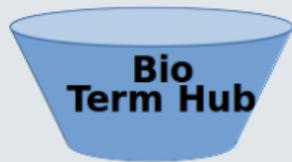
References



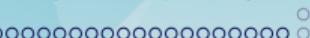
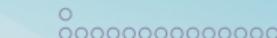
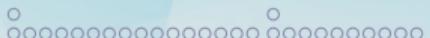
Cell Ontology



Gene Ontology



Bio Term Hub



Resource Selection

Please select the resources to be included:

[select all](#)

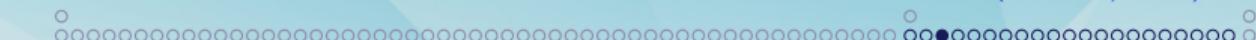
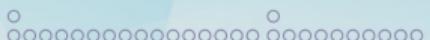
<input type="checkbox"/> Cell Ontology	(→ source)	Update available.	<input type="button" value="Update"/>
<input type="checkbox"/> Cellosaurus	(→ source)	Up-to-date.	<input type="button" value="..."/>
<input type="checkbox"/> ChEBI	(→ source)	Up-to-date.	<input type="button" value="..."/>
<input type="checkbox"/> CTD chemicals	(→ source)	Up-to-date.	<input type="button" value="..."/>
<input type="checkbox"/> CTD diseases	(→ source)	Up-to-date.	<input type="button" value="..."/>
<input type="checkbox"/> EntrezGene	(→ source)	Update available.	<input type="button" value="Update"/>
<input type="checkbox"/> Gene Ontology	(→ source)	Update available.	<input type="button" value="Update"/>
<input type="checkbox"/> MeSH	(→ source)	Up-to-date.	<input type="button" value="..."/>
<input type="checkbox"/> NCBI Taxonomy	(→ source)	Update available.	<input type="button" value="Update"/>
<input type="checkbox"/> Protein Ontology	(→ source)	Up-to-date.	<input type="button" value="..."/>
<input type="checkbox"/> Sequence Ontology	(→ source)	Up-to-date.	<input type="button" value="..."/>
<input type="checkbox"/> Swiss-Prot	(→ source)	Up-to-date.	<input type="button" value="..."/>
<input type="checkbox"/> Uberon	(→ source)	Up-to-date.	<input type="button" value="..."/>

The Bio Term Hub provides access to biomedical terminology resources in a unified format. We envision its main use as a basis for text mining systems.

Select any of the external resources on the left to obtain a custom terminology list. The list is a tab-separated table including terms (names, synonyms), preferred name, concept identifier, entity type, and original resource. The contents of the last two fields can be modified using replacement patterns (uncover the **Renaming** section below using the triangle button). Submit your request with the **Create resource** button at the bottom of the page. Afterwards you will be offered the opportunity to annotate a text with the selected resources.

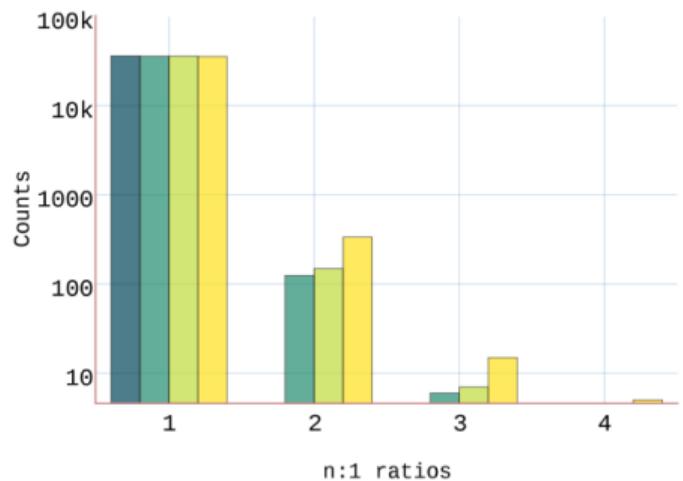
Aggregating, filtering, reformatting, and optionally renaming all this information takes time. Creation time mainly depends on the size of each resource. Compiling a list of all terminologies except for EntrezGene takes approximately 90 seconds; however, EntrezGene is very large, and processing it requires 3 to 4 minutes.

The requested terminology list is compiled on the fly, based on local copies of the external terminology resources. Whenever one of the local copies becomes out-of-date with respect to its remote source, an **Update** button is shown next to the corresponding resource name. This button triggers downloading the latest version to our server. Please note that, due to data volume and bandwidth restrictions, updating may take several minutes for some of the resources.

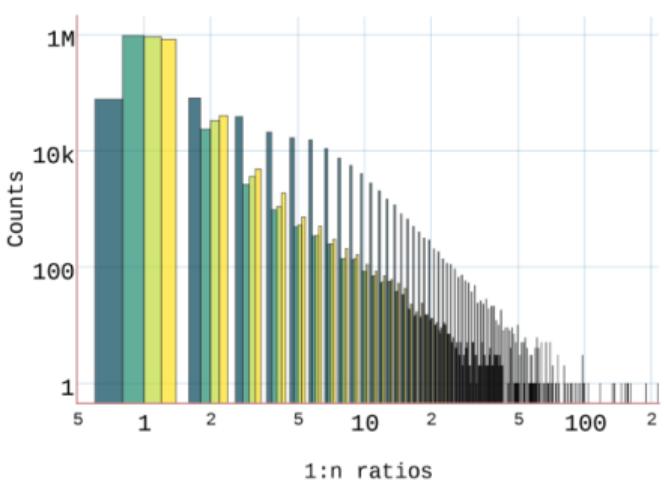


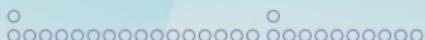
Bio Term Hub

Cell lines



Chemicals

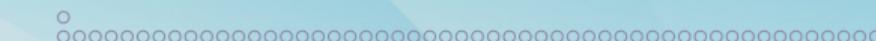
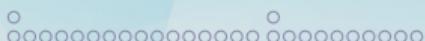




BTH: Term Statistics

Tabelle: Overview of Termfile Statistics

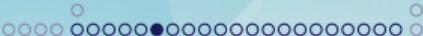
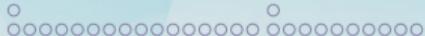
	genes/proteins	chemicals	diseases	species	cell lines	all
Terms in Resource	10,429,162	979,418	67,614	1,333,903	36,249	12,846,346
Avg. term length (number letters)	11.73	37.49	26.98	22.87	7.611	14.92
Avg. terms/ID	1.1455	3.545	6.018	1.326	1.000	1.328
Avg. IDs/term	1.371	1.049	1.000	1.003	1.004	1.306
Avg. IDs/term (case insensitive)	1.383	1.062	1.000	1.003	1.004	1.317
Avg. IDs/term (case insensitive, stripped)	1.387	1.086	1.000	1.006	1.010	1.324



BTB: Term confusion matrix

Overlap among entity types

	chemical	sequence	gene/ protein	cellular component	cell	cell line	molecular function	biological process	organism	disease
chemical		7.46%	27.19%	4.51%	0.26%	2.56%	3.70%	0.08%	0.76%	1.61%
sequence	23.68%		23.72%	15.41%	0.00%	1.69%	4.18%	0.29%	1.10%	0.36%
gene/protein	8.89%	2.44%		2.77%	1.12%	3.21%	2.22%	0.12%	2.29%	1.77%
cellular component	29.63%	31.84%	55.70%		16.74%	1.60%	4.43%	0.60%	0.22%	1.27%
cell	3.97%	0.00%	51.83%	38.57%		4.38%	0.82%	0.45%	0.06%	3.49%
cell line	17.75%	3.69%	67.86%	1.69%	2.00%		2.10%	0.19%	3.07%	5.77%
molecular function	34.50%	12.27%	63.34%	6.29%	0.51%	2.82%		0.68%	1.09%	2.30%
biological process	0.41%	0.49%	2.07%	0.50%	0.16%	0.15%	0.40%		0.03%	3.79%
organism	3.14%	1.42%	28.79%	0.14%	0.02%	1.83%	0.48%	0.02%		0.35%
disease	4.09%	0.29%	13.75%	0.49%	0.59%	2.12%	0.63%	1.78%	0.21%	



OGER

Human prostate cancer metastases target the hematopoietic stem cell niche to establish footholds in mouse bone marrow.

HSC homing, quiescence, and self-renewal depend on the bone marrow HSC niche.

A large proportion of solid tumor metastases are bone metastases, known to usurp HSC homing pathways to establish footholds in the bone marrow.

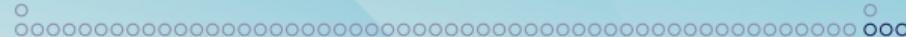
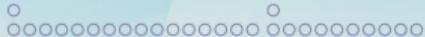
However, it is not clear whether tumors target the HSC niche during metastasis.

Legend

- disease
- chemical
- sequence
- gene/protein
- biological_process
- organism
- cell

<https://covid19.nlp.idsia.ch/oger-rest.html>

<https://oger.nlp.idsia.ch/>

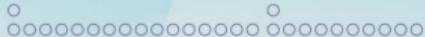


OGER: annotation service

The OntoGenes Biomedical Entity Recogniser (OGER)

- RESTful web service, using BTH terminologies
- Allows annotation of a collection of documents.
- Evaluated in the Bio Text Mining services challenge BioCreative/TIPS
 - best results according to several of the evaluation metrics.

<https://covid19.nlp.idsia.ch/oger-rest.html>

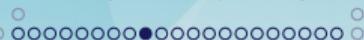
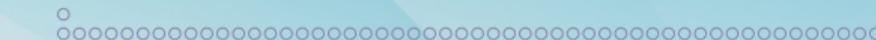


OGER: annotation service

- Annotates input text with entities from the BTH
 - Except EntrezGene
- Can be used as a web demo (for annotation of single articles) or as a web service (batch).
- Input: PubMed, PubMed Central, Free Text
- Formats: text(I), BioC (I/O), pxml (I), tsv (O), brat (O), odin-xml (O)

<https://oger.nlp.idsia.ch/>

Note: user-provided terminologies can be used, but this is not yet supported by the interface and web service.



BioCreative V.5 / TIPS



seconds/document

1.06943 s

1st

seconds/Byte

0.00086 s

1st

seconds/annotation

0.07227 s



annotations/document

14.7923



time between failures

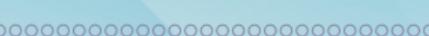
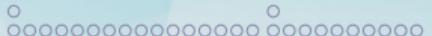
[no failure]

1st (shared)

time to repair

[no failure]

1st (shared)



OntoGene in BioCreative II.5

Fabio Rinaldi, Gerold Schneider, Kaarel Kaljurand, Simon Clematide, Thérèse Vachon, and Martin Remacker

Abstract—We describe a system for the detection of mentions of protein-protein interactions in the biomedical scientific literature. The original system was developed as a part of the OntoGene project, which focuses on using advanced computational linguistic techniques for text mining applications in the biomedical domain. In this paper, we focus in particular on the participation to the BioCreative II.5 challenge, where the OntoGene system achieved best-ranked results. Additionally, we describe a feature-analysis experiment performed after the challenge, which shows the unexpected result that one single feature alone performs better than the combination of features used in the challenge.

Index Terms—Biomedical text mining, Natural Language Processing (NLP), protein interactions, BioCreative.

1 INTRODUCTION

A s a way to cope with the constantly increasing generation of biomedical knowledge in biology, some organizations maintain various types of databases that aim at collecting the most significant information in a specific area. For example, UniProt/SwissProt [1] collects information on all known proteins. MINT[2] and IntAct[3] are the databases collecting protein interactions. Most of the information in these databases is derived from the primary literature by a process of manual revision known as “literature curation.”

Typically, a text mining system is designed to locate relevant entities mentioned in the literature (e.g., by PubMed database of biomedical literature), not only by finding relevant articles (which is the task of information retrieval systems), but also by extracting very specific information of interest to the user, such as, for example, protein-protein interactions (PPI)¹. Text mining systems are part of a variety of different approaches, ranging from simple baseline methods to more complex natural language approaches, with varying degrees of success.

The work presented here is part of a larger effort undertaken in the OntoGene project [6] aimed at improving biomedical text mining through the usage of advanced natural language processing techniques. The results of the entity detection feed directly into the process of identification of protein interactions. Our approach relies upon information delivered by a pipeline of NLP tools, including sentence splitting, tokenization, part of speech tagging, term

recognition, noun and verb phrase chunking, and a dependency-based syntactic analysis of input text [7]. There are specific [9] tasks into document constituent boundaries defined by previously identified multiword entities. Therefore, the richness of the entity annotation has a direct beneficial impact on the performance of the parser, and thus, leads to better recognition of interactions [10].

There have recently been numerous results showing the potential of dependency-based language analysis for text mining. Fysall et al. [11] describe a manually annotated corpus, which includes a dependency-based analysis of each sentence. Clegg and Sproat [12] use dependency graphs in order to build a mark's four publicly available natural language processors. Fundel et al. [13] describe a large-scale relation mining application based upon the Stanford Lexicalized Parser.

In this paper, we first describe in Section 2 the overall architecture of our text mining system. We then present in Section 3 the results obtained in the recent BioCreative II.5 competition. Additionally, in Section 4, we briefly sketch how the system could be used in the curation process.

2 THE ONTOGENE TEXT MINING SYSTEM

In this section, we provide an overall description of the OntoGene text mining environment, with a specific focus on its application to the detection of PPI. In Section 2.1, we explain the process used to automatically annotate different types of entities, and group them to reference identifiers (IDs). Section 2.2 illustrates how to extract information about the focus organisms mentioned in the articles and how we use it to disambiguate protein mentions. In Section 2.3, we describe our approach to the detection of interactions among entities (proteins in particular).

2.1 Detection and Grounding of Domain Entities

In this section, we describe our approach to the problem of detecting names of relevant domain entities in biomedical literature (we consider, in particular, proteins, genes, species, experimental methods, and cell lines) and grounding them to widely accepted IDs assigned by four different knowledge bases: UniProt Knowledgebase [14], National

¹ Through surveys of the role of Text Mining systems in biology can be found in [4] and [5].

² F. Rinaldi, G. Schneider, K. Kaljurand, and S. Clematide are with the Institute of Computational Linguistics, University of Zurich, Rosenbergstrasse 34, CH-8050 Zurich, Switzerland.
E-mail: fabio.rinaldi@inf.usz.ch; gerold.schneider@inf.usz.ch; kaarel.kaljurand@inf.usz.ch; simon.clematide@inf.usz.ch.

³ T. Vachon and M. Remacker are with Novartis Pharma AG, NIBB/NIAID/NIH Mining Seminar, CH-4049 Basel, Switzerland.
E-mail: therese.vachon@novartis.com; martin.remacker@novartis.com.

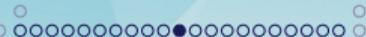
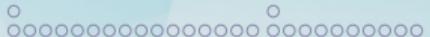
Manuscript received 11 Jun. 2009; revised 2 Apr. 2010; accepted 3 May 2010; published online 27 May 2010.
For information on obtaining reprints of this article, please send e-mail to: tcv@nibb.it.

Digital Object Identifier no. 10.1109/TCBB.2010.201509.

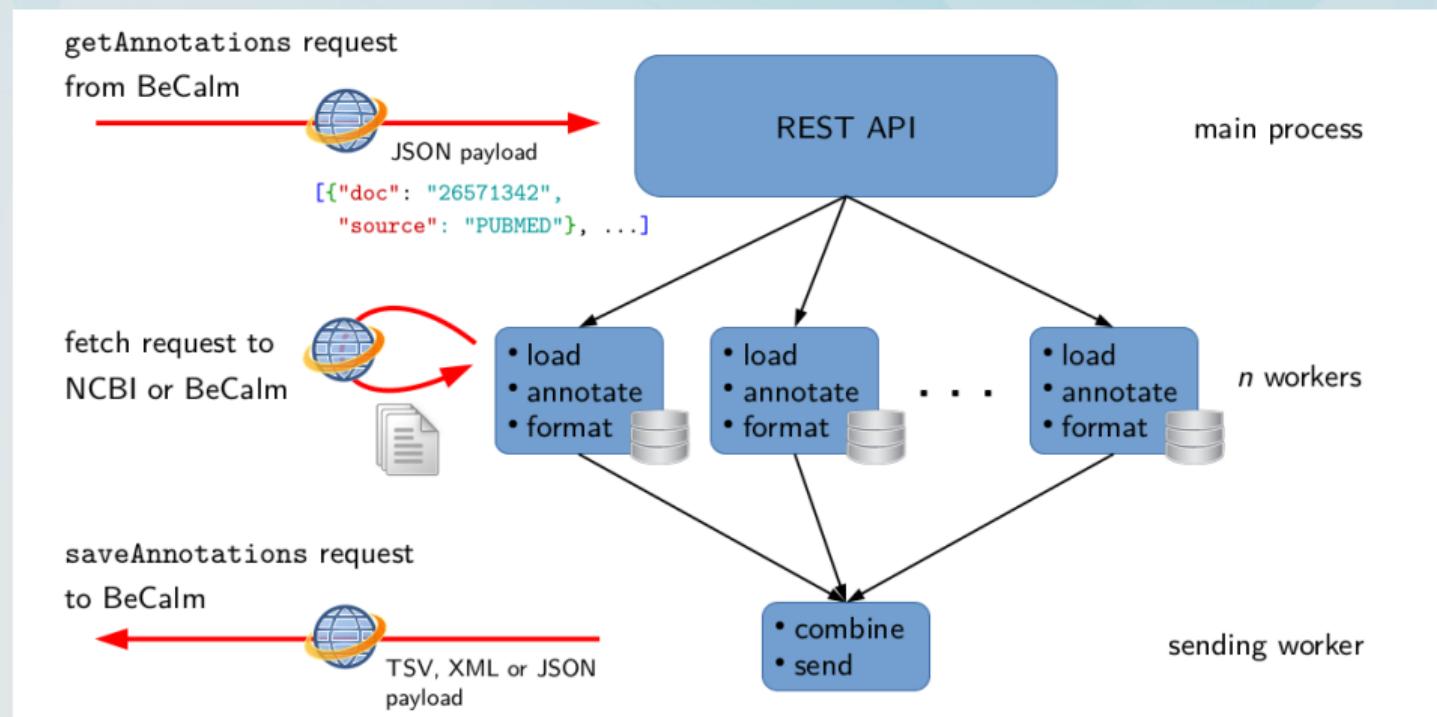
Previous history . . .

- [2006] BioCreative II: PPI (3rd), IMT (best)
- [2009] BioCreative II.5 PPI (best results); BioNLP
- [2010] BioCreative III: ACT, IMT, IAT
- [2011] CALBC (large scale entity extraction), BioNLP
- [2012] CTD task at BioCreative 2012
- [2013] BioCreative IV: BioC, CTD, IAT

<http://www.biomext.org/>



OGER in TIPS



NOTE: Just won an OpenMinTeD tender! (17'500 EUR)



Disambiguation challenges

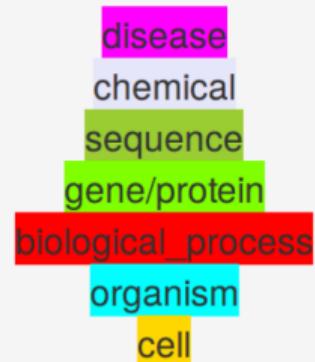
Human prostate cancer metastases target the hematopoietic stem cell niche to establish footholds in mouse bone marrow.

HSC homing, quiescence, and self-renewal depend on the bone marrow HSC niche.

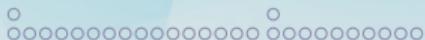
A large proportion of solid tumor metastases are bone metastases, known to usurp HSC homing pathways to establish footholds in the bone marrow.

However, it is not clear whether tumors target the HSC niche during metastasis.

Legend

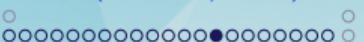
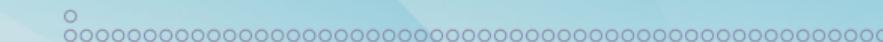


gene/protein	galactoside 2-alpha-L-fucosyltransferase 1	PR:000007702	Protein Ontology
cellular component	Hedgehog signaling complex	GO:0035301	Gene Ontology
gene/protein	chaperone protein HscA	PR:000022925	Protein Ontology
chemical	N-(3-carboxypropanoyl)-N-hydroxycadaverine	CHEBI:50443	ChEBI



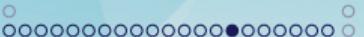
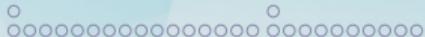
Term ambiguity

Name/Gene ID	Description	Location	Aliases
<input type="checkbox"/> thiS ID: 2847702	immediate sulfur donor in thiazole formation [<i>Escherichia coli</i> str. K-12 substr. MG1655]	NC_000913.3 (4192637..4192837, complement)	b4407, ECK3983, JW3955
<input type="checkbox"/> thiS ID: 939811	sulfur carrier protein ThiS [<i>Bacillus subtilis</i> subsp. <i>subtilis</i> str. 168]	NC_000964.3 (1244844..1245044)	BSU11680
<input type="checkbox"/> thiS ID: 4524717	ThiS [<i>Thalassiosira pseudonana</i>]	NC_008589.1 (84740..84952)	ThpsCp091, ycf40
<input type="checkbox"/> thiS ID: 4524549	ThiS [<i>Phaeodactylum tricornutum</i>]	NC_008588.1 (81332..81544)	PhtrCp091
<input type="checkbox"/> This ID: 33350719	This [<i>Sheathia arcuata</i>]	NC_035231.1 (47491..47706, complement)	CGW65_pgp144



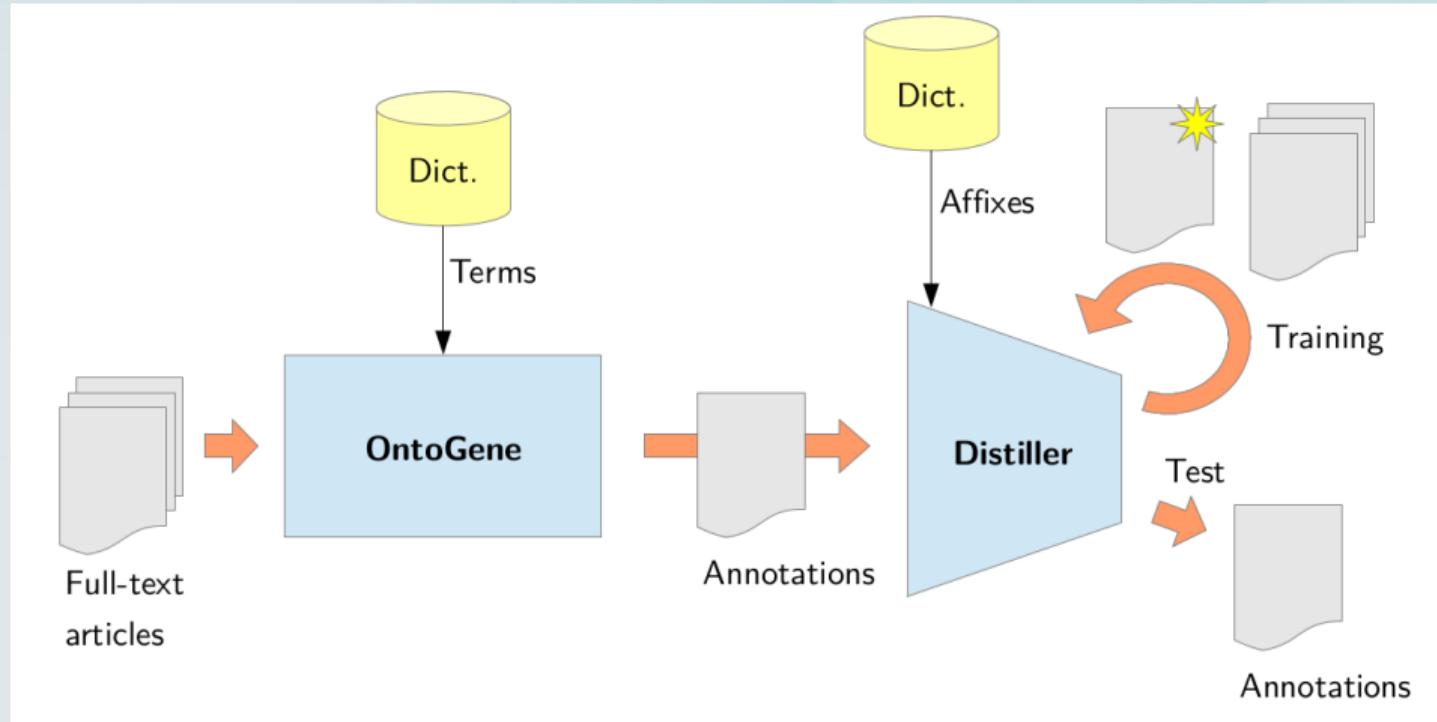
Term ambiguity

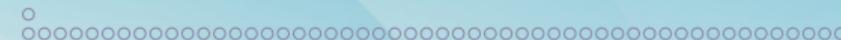
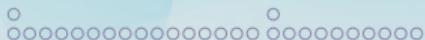
Name/Gene ID	Description	Location	Aliases	MIM
<input type="checkbox"/> WAS ID: 7454	Wiskott-Aldrich syndrome [<i>Homo sapiens</i> (human)]	Chromosome X, NC_000023.11 (48683753..48691427)	IMD2, SCNX, THC, THC1P, WASPA, WAS	300392
<input type="checkbox"/> Was ID: 22376	Wiskott-Aldrich syndrome [<i>Mus musculus</i> (house mouse)]	Chromosome X, NC_000086.7 (8081466..8090491, complement)	U42471p, Was	
<input type="checkbox"/> Was ID: 317371	Wiskott-Aldrich syndrome [<i>Rattus norvegicus</i> (Norway rat)]	Chromosome X, NC_005120.4 (15155246..15164099)	WASP	
<input type="checkbox"/> WAS ID: 107131173	Wiskott-Aldrich syndrome [<i>Bos taurus</i> (cattle)]	Chromosome X, AC_000187.1 (91842853..91848367)		



Term ambiguity

Name/Gene ID	Description	Location	Aliases	MIM
<input type="checkbox"/> CAT ID: 847	catalase [<i>Homo sapiens</i> (human)]	Chromosome 11, NC_000011.10 (34438925..34472060)		115500
<input type="checkbox"/> Cat ID: 40048	Catalase [<i>Drosophila melanogaster</i> (fruit fly)]	Chromosome 3L, NT_037436.4 (18822604..18828188)	Dmel(CG6871, CAT, CATA, CG6871, CT21282A, DMCATHPO, DROCATHPO, Dmel\CG6871, U00145, bs36h11.y1, cat, catalase, Cat	
<input type="checkbox"/> CAT ID: 101093891	catalase [<i>Felis catus</i> (domestic cat)]	Chromosome D1, NC_018732.2 (88501679..88538456)		
<input type="checkbox"/> CAT ID: 531682	catalase [<i>Bos taurus</i> (cattle)]	Chromosome 15, AC_000172.1		





CRAFT corpus

Corpus

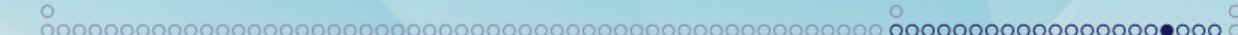
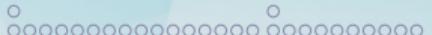
- 67 full-text articles (PubMed Central)
- 100 k manually verified concept annotations, completed in 2012
- refers to 7 ontologies/terminologies:
ChEBI, CL, Entrez Gene, GO, NCBI taxonomy, PRO, SO

Usage

CRAFT is being used to date^{1,2} for evaluating concept recognition systems.

¹Christopher Funk et al. (2014). "Large-scale biomedical concept recognition: an evaluation of current automatic annotators and their parameters". In: *BMC Bioinformatics* 15.1, pp. 1–29

²Eugene Tseytlin et al. (2016). "NOBLE – Flexible concept recognition for large-scale biomedical text mining". In: *BMCG Bioinformatics* 17.1, pp. 1–15



Bastida et al. Journal of Biomedical Semantics (2015) 6:11
DOI 10.1186/s13396-015-0074-z

Journal of
Biomedical Semantics

RESEARCH

Open Access



Entity recognition in the biomedical domain using a hybrid approach

Marco Bastida^{1,2}, Lenz Furrer^{3,4}, Carlo Tassio¹ and Fabio Rinaldi^{1,2*}

Abstract

Background: This article describes a high-recall, high-precision approach for the extraction of biomedical entities from scientific articles.

Method: The approach uses a two-stage pipeline, combining a dictionary-based entity recognizer with a machine-learning classifier. First, the OGER entity recognizer, which has a bias towards high recall, annotates the terms that appear in selected domain ontologies. Subsequently, the Distiller framework uses this information as a feature for a machine learning algorithm to select the relevant entities only. For this step, we compare two different supervised machine-learning algorithms, Conditional Random Fields and Neural Networks.

Results: In an independent evaluation on the CHASE corpus, we test the performance of the combined systems. After the first stage of OGER, the recall of the system is 90%. Subsequently, the second stage, involving an ensemble of classifiers and logistic regression, increases the recall to 96%. Our final system achieves a 98% F-score, using a Neural Network-based filtering. It achieves an overall precision of 98% at a recall of 68% on the named entity recognition task, and a precision of 91% at a recall of 69% on the concept recognition task.

Conclusion: These results are to our knowledge the best reported to far in this particular task.

Keywords: Name/entity recognition, Text mining, Machine learning, Natural language processing

Background

The scientific community in the biomedical domain is a vibrant community, producing a large amount of scientific findings in the form of data, publications, reports, and so on, each year, making it difficult for scholars to find the right information in this large sea of knowledge.

To tackle this problem, researchers have developed different text mining techniques with the goal of detecting the relevant entities in the biomedical corpus. The main focus is in the technique called Named Entity Recognition (henceforth NER), which solves the problem of detecting terms belonging to a limited set of predefined entity types.

NER can be performed on both "generic" documents, to recognize concepts like person, state or location, or on technical documents, to recognize concepts like cells.

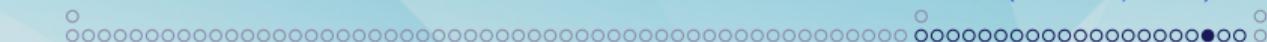
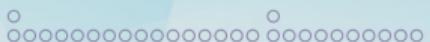
*Correspondence: rinaldi@inf.unibz.it

^{1,2}University of Zurich, Institute of Computational Linguistics and Swiss Institute of Bioinformatics, Winterthurerstrasse 190, 8057 Zurich, Switzerland
Full list of author information is available at the end of the article



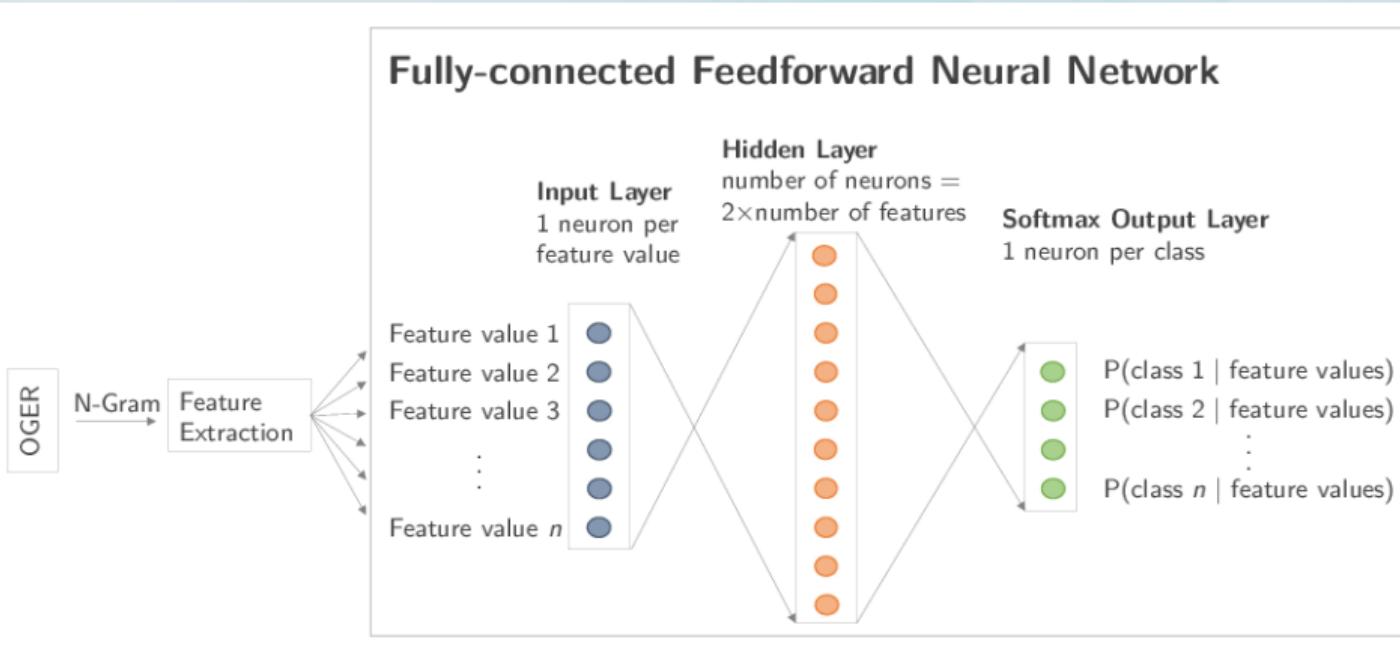
© The Author(s). 2015. Open Access This article is distributed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated.

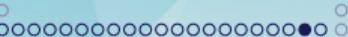
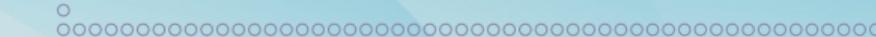
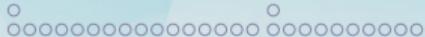
	Neural Network	Conditional Random Fields
Implementation		
Software	R [47], nnet library	CRFSuite [48]
Model parameters	1 hidden layer of size $2 \times (\text{number of input features})$, softmax output layer	Training algorithm: averaged perceptron, default epsilon, 2 words window
Input	n-grams selected by OGER	single tokens
Features		
Candidate character count	count	—
Candidate is all uppercase	label yes/no	label yes/no
Candidate is all lowercase	label yes/no	label yes/no
Candidate contains Greek (i.e. "alpha", α)	label yes/no	label yes/no
Candidate contains dashes ('-')	count	label yes/no
Candidate contains numbers	count	label yes/no
Candidate ends with a number	label yes/no	label yes/no
Candidate contains capital letter not in first position	label yes/no	label yes/no
Candidate contains lowercase characters	count	label yes/no
Candidate contains uppercase characters	count	label yes/no
Candidate contains spaces	count	label yes/no
Candidate contains symbols	count	label yes/no
2-3 character affixes appearing in an ontology in [36]	normalized frequency	label yes/no
Candidate is symbol	—	label yes/no
Candidate's part-of-speech	—	yes, using [49]
Candidate's stem	—	yes, using [50]
Candidate pre-selected by OGER	—	yes (see Section 2.5)
Total features	36	about 2.8 million
Tagging speed	1286 tokens/sec	632 tokens/sec



New JointNN filter

Fully-connected Feedforward Neural Network



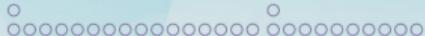


Results

Table 4: F1-scores

Entity Type	OGER	OGER+DistillerNN	OGER+JointNN
All	0.51	0.70	0.71
Chemicals	0.55	0.77	0.79
Cells	0.80	0.76	0.72
Biological processes/molecular functions	0.30	0.35	0.45
Cellular components	0.55	0.70	0.73
Organisms	0.44	0.94	0.95
Proteins	0.62	0.80	0.79
Sequences	0.54	0.75	0.73

[Anna Jancso, Lenz Furrer, Fabio Rinaldi, in preparation]

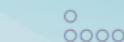


Conclusions

- Biomedical Text Mining is indispensable to make sense of huge amounts of text: e.g. scientific literature, clinical records, social media
- The complexity of human language makes it difficult for machines to extract meaning out of text
- The biomedical field has been at the forefront of research in text mining, with numerous competitions and shared resources
- Recent advances in NLP, based on deep learning, have generated a quantum leap in our ability to extract information from text

Introduction

The Biomedical Literature Methods



Our Tools ("OntoGene / IDSIA")

References



Topic

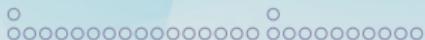
Introduction

The Biomedical Literature

Methods

Our Tools ("OntoGene / IDSIA")

References



References

-  Ashburner, M., Ball, C., Blake, J., Botstein, D., Butler, H., Cherry, J., Davis, A., Dolinski, K., Dwight, S., Eppig, J., Harris, M., Hill, D., Issel-Tarver, L., Kasarskis, A., Lewis, S., Matese, J., Richardson, J., Ringwald, M., Rubin, G., and Sherlock, G. (2000). Gene ontology: tool for the unification of biology. The Gene Ontology Consortium. *Nat. Genet.*, 25:25–29.
-  Chang, J. T., Schtze, H., and Altman, R. B. (2002). Creating an online dictionary of abbreviations from MEDLINE. *J Am Med Inform Assoc*, 9(6):612–20. PMID:12386112.
-  Chen, H. and Sharp, B. M. (2004). Content-rich biological network constructed by mining PubMed abstracts. *BMC Bioinformatics*, 5:147. PMID:15473905.