# Towards Biomedical Neurosymbolic AI: From Semantic Knowledge Infrastructure to Explainable Predictions

# FAIR

**Michel Dumontier, PhD**

Distinguished Professor of Data Science
Founder and Director, Institute of Data Science
Department of Advanced Computing Sciences
Maastricht University

EGC: Strasbourg :: 29-01-2025

# Drug discovery

**Goal is to discover a molecule that can efficiently correct/modify an abnormal/undesirable trait associated with minimal adverse effects.**
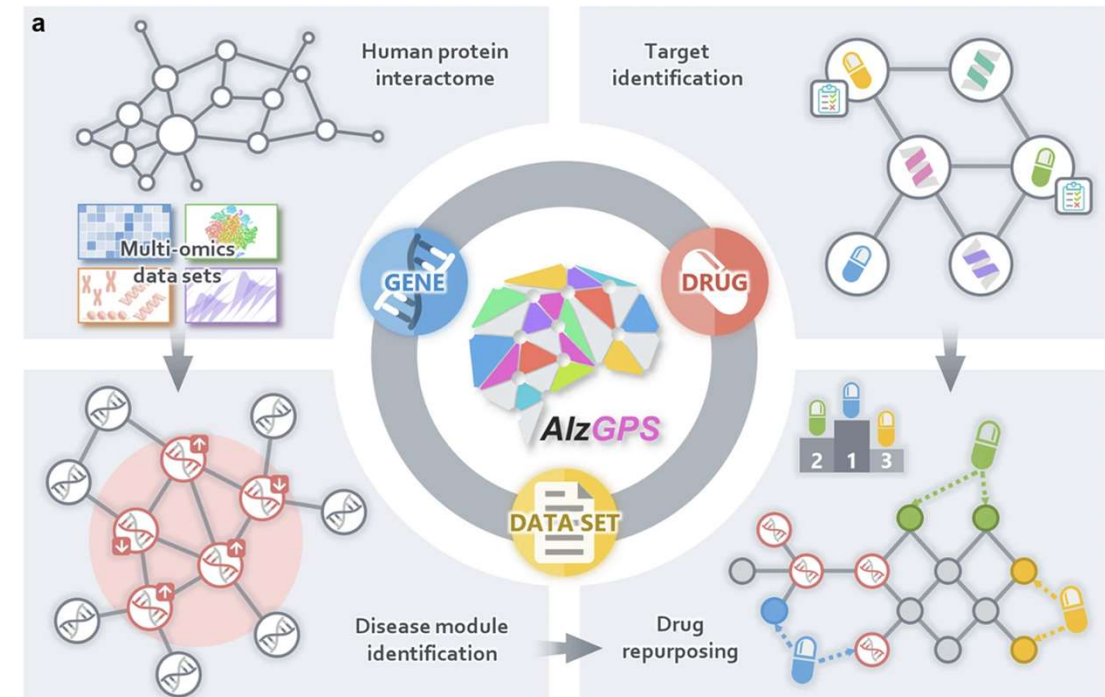
22,000+ human diseases*, <1000 with approved treatments. Treatments are not always effective.

Much more work is needed to find effective therapeutics across the space of all diseases

Drug discovery and development takes years and lots of $$$. Involves:
- studying the disease mechanism
- identifying candidate targets for intervention
- developing/testing molecules
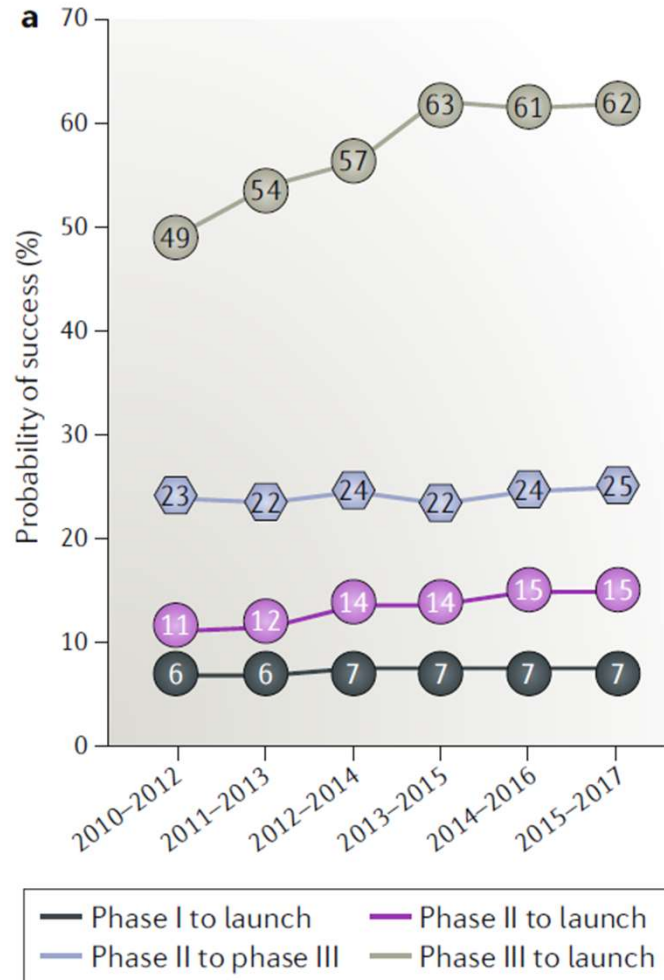- performing in vitro/cell/tissue/animal studies
- human validation



*Alz Res Therapy* **13**, 24 (2021).

# THE CLINICAL-TRIAL CLIFF

Drug companies are removing more compounds from the pipeline at all levels of testing than ever before.

Probability to launch

**To what extent will AI improve the success of developing new treatments?**

# A new generation of companies ( and researchers ) are betting on data and AI

***significant effort*** is needed to **find** the right data, make **sense** of them, and **use** them for a new purpose

# Data scientists could be more productive



## What data scientists spend the most time doing

- Building training sets: 3%
- Cleaning and organizing data: 60%
- Collecting data sets: 19%
- Mining data for patterns: 9%
- Refining algorithms: 4%
- Other: 5%

**data** remain challenging to access and reuse:

- difficult to obtain
- poorly described
- in different formats
- hard to integrate with other data

**AI** have significant limitations:

- built from limited, biased, or un-representational data
- aren't robust to different inputs
- have a hard time predicting out of distribution learning
- may not be able to explain or justify outputs

**Translational Failure**

# Human Machine collaboration
# is crucial to our future work

# Machines
need to be able to discover and reuse data
(and arguably any digital resource)

high quality, machine accessible, linked, (meta)data from multiple sources and data types

Trustworthy, data-oriented, explainable AI models and predictions

Translational Success

# Research Directions

The totality of (digitized) biomedical knowledge and analytics to:

**i) answer questions** about what we know and what we don't know (*but should*)

**ii) create robust models to predict, explain and justify biomedical phenomena**

**iii) tools for human-AI collaboration** to create, maintain, correct, and complete knowledge

**Knowledge Infrastructure** ⟷ **Explainable Predictions**

*FAIR*
*Data & Services*

*Neurosymbolic AI*

# Findable    Accessible    Interoperable    Reusable

# The FAIR Guiding Principles for scientific data management and stewardship

Mark D. Wilkinson, Michel Dumontier, … Barend Mons ✉   + Show authors

## Box 2 | The FAIR Guiding Principles

**To be Findable:**
F1. (meta)data are assigned a globally unique and persistent identifier
F2. data are described with rich metadata (defined by R1 below)
F3. metadata clearly and explicitly include the identifier of the data it describes
F4. (meta)data are registered or indexed in a searchable resource

**To be Accessible:**
A1. (meta)data are retrievable by their identifier using a standardized communications protocol
A1.1 the protocol is open, free, and universally implementable
A1.2 the protocol allows for an authentication and authorization procedure, where necessary
A2. metadata are accessible, even when the data are no longer available

**To be Interoperable:**
I1. (meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation.
I2. (meta)data use vocabularies that follow FAIR principles
I3. (meta)data include qualified references to other (meta)data

**To be Reusable:**
R1. meta(data) are richly described with a plurality of accurate and relevant attributes
R1.1. (meta)data are released with a clear and accessible data usage license
R1.2. (meta)data are associated with detailed provenance
R1.3. (meta)data meet domain-relevant community standards

EUROPEAN COMMISSION
Press Release Database
European Commission > Press releases database > Press Release details
**European Commission - Statement**
**G20 Leaders' Communique Hangzhou Summit**
Hangzhou, 5 September 2016
1. We, the Leaders of the G20, met in Hangzhou, China on 4-5 September 2016.

G7 2017 ITALIA

**Annex 4:**
**G7 Expert Group on Open Science**
Turin, Italy, September 28, 2017

Realising the European Open Science Cloud
First report and recommendations of the Commission High Level Expert Group on the European Open Science Cloud

Final Report and Action Plan from the European Commission Expert Group on FAIR Data

**TURNING FAIR INTO REALITY**
2018

http://www.nature.com/articles/sdata201618

15

# Making FAIR Data

**1. Collect**

Data

**2. Describe**

Standardized Metadata

use standard metadata format

use ontologies + vocabularies

add provenance, license for data + metadata

**3. Transform**

Standardized Data

Use standard data format

use ontologies + vocabularies

**4. Publish**

Findable Accessible Interoperable Reusable

Data Repository

Persistent Metadata Identifier

Standardized Metadata

Persistent Data Identifier

Standardized Data

# Communities are publishing recipes
# to make FAIR data

# How do we know it's FAIR?

- **FAIR Enough** is a system to perform automated assessment of the technical quality of the FAIRness implementation.
- Uses a collections of metrics, implemented as web services.
- Fast owing to parallel execution
- Keeps track of past assessments to monitor status
- Offers search and query services
- Anybody can extend via service based framework
- Open source and Docker deployable



https://fair-enough.semanticscience.org

# DataverseNL

(DANS)

# Welcome to DataverseNL

Store, share and publish research data online. Use the slider below to access the dataverses of the DataverseNL partners.

Please use the DataverseNL demo for demonstration, training and testing purposes.

If you have questions about how to use DataverseNL, please contact your institution directly via one of the email addresses on this page.

‹

Rijksdienst voor het Cultureel Erfgoed
Ministerie van Onderwijs, Cultuur en Wetenschap

**Rijksdienst voor het Cultureel Erfgoed**

Utrecht University

**Utrecht University**

Maastricht University

**Maastricht University**

rijksuniversiteit groningen

**University of Groningen**

›

Search this dataverse...    🔍    Advanced Search

☑ 👥 **Dataverses (657)**
☑ 📄 **Datasets (7,190)**
☐ 📄 Files (189,511)

**Dataverse Category**
Research Project (170)
Organization or Institution (149)
Department (148)
Research Group (77)
Laboratory (18)

More...

**Publication Year**
2024 (358)
2023 (654)
2022 (749)
2021 (4,137)
2020 (372)

More...

1 to 10 of 7,847 Results          ↕ Sort ▾

**Dendrochronologisch onderzoek Groene Fietsgordel – inrichtingsdossier 3 Ver-Assebroek en Magdalenapad (BAAC project 23.1060)**
Mar 20, 2024 - BAAC

Doeve, Petra, 2024, "Dendrochronologisch onderzoek Groene Fietsgordel – inrichtingsdossier 3 Ver-Assebroek en Magdalenapad (BAAC project 23.1060)", https://doi.org/10.34894/UQS0FK, DataverseNL, V2

Twee houtmonsters (M1 en M2) van mogelijke brugpijlers zijn ingezet voor dendrochronologisch onderzoek, met als doel de ouderdom van het hout en het kapjaar van de boom zo nauwkeurig mogelijk te bepalen. De houten paal zijn aangetroffen in een archeologisch beschermde zone in Ver...

**Results of semi-systematic literature review on the application of Q-methodology in the rural Global South, as part of the DARE-TU PhD project**
Mar 19, 2024 - ITP project

Intriago Zambrano, Juan Carlo, 2024, "Results of semi-systematic literature review on the application of Q-methodology in the rural Global South, as part of the DARE-TU PhD project", https://doi.org/10.34894/K252KB, DataverseNL, V2

This dataset of the DARE-TU PhD project comprises the results of the systematic literature review of 47 selected studies on the application of Q methodology in the rural Global South. It encompasses a number of qualitative and quantitative information on the experimental design,...

Tasks   Sizes   Sub-tasks   Languages   Licenses   Other

🔍 Filter Tasks by name

**Multimodal**

⊞ Feature Extraction        📝 Text-to-Image

📝 Image-to-Text            ⊡ Text-to-Video

⊡ Visual Question Answering    ⧉ Graph Machine Learning

**Computer Vision**

⬡ Depth Estimation         ⊡ Image Classification

⊡ Object Detection         ⊠ Image Segmentation

⊡ Image-to-Image           ⊡ Unconditional Image Generation

⊡ Video Classification     ⊡ Zero-Shot Image Classification

**Natural Language Processing**

⊡ Text Classification      ⊡ Token Classification

⊞ Table Question Answering    ⊡ Question Answering

✳ Zero-Shot Classification    ⊠ Translation

⊡ Summarization            ⊡ Conversational

⊡ Text Generation          ⊡ Text2Text Generation

⊡ Fill-Mask    ⊡ Sentence Similarity    ⊡ Table to Text

⊡ Multiple Choice          ⊡ Text Retrieval

## Datasets  71,980

🔍 Filter by name

▤ **open-web-math/open-web-math**
⊞ Viewer • Updated 5 days ago • ⬇ 576 • ♡ 149

▤ **openbmb/UltraFeedback**
⊞ Viewer • Updated 23 days ago • ⬇ 856 • ♡ 111

▤ **EleutherAI/proof-pile-2**
⊞ Viewer • Updated 5 days ago • ⬇ 163 • ♡ 37

▤ **laion/dalle-3-dataset**
⊞ Viewer • Updated about 3 hours ago • ⬇ 1.09k • ♡ 141

▤ **THUDM/AgentInstruct**
⊞ Viewer • Updated 3 days ago • ⬇ 61 • ♡ 21

▤ **ehartford/dolphin**
◉ Preview • Updated 28 days ago • ⬇ 472 • ♡ 211

▤ **akjindal53244/Arithmo-Data**
◉ Preview • Updated about 12 hours ago • ⬇ 21 • ♡ 13

▤ **OpenAssistant/oasst1**
⊞ Viewer • Updated May 2 • ⬇ 9k • ♡ 1.05k

---

🗄 Datasets: 🐺 laion / **dalle-3-dataset** ⧉   ♡ like  141

Languages: 🌐 English   Tags: image-text-dataset   synthetic-dataset   License: 🏛 cc0-1.0

🗂 **Dataset card**    ⊟ Files and versions    🤗 Community 6

⊞ **Dataset Viewer**        ↻ Auto-converted to Parquet  </> API  ⊞ Go to dataset viewer

Split

train (4.71k rows)                                              ⌄

| caption string | image image | link string |
|---|---|---|
| Amidst the towering... |  | https://cdn.discordapp.com/attachments/1158354590463447092/1163229131371192431/9e18e126-ec8e-4c62-b512-hm=ec95543d046be2f1879c569890a298a043070d34eebcfa3106d6141e24d0133f& |
| Amidst the towering... |  | https://cdn.discordapp.com/attachments/1158354590463447092/1163229131752865832/27f8b6c9-e7ca-47fd-a4ea-hm=f6ff694fd0c202d94ff0a4bf2200caf31c1cd1a909d3922981d114a107c97f1c& |
| Amidst the towering... |  | https://cdn.discordapp.com/attachments/1158354590463447092/1163229132277161984/ed3fbabf-16d1-471e-9c33-hm=d9ac76647142a1add3aba2beb6e112051c727b041d3224092873bd1ead215647& |
| Amidst the towering... |  | https://cdn.discordapp.com/attachments/1158354590463447092/1163229132583338104/265825c4-5e29-4ca8-81de-hm=c40571b29059e389ea9db2d5c44785779e0c93c4697e2e5c0f05bb8e9f275688& |
| Wide anime- |  | https://cdn.discordapp.com/attachments/1158354590463447092/1163228854257733782/3fd3de66-0b14-43d5-ba92- |

◁ Previous  **1**  2  3  ...  48  Next ▷

## Dataset Card for LAION DALL·E 3 Discord Dataset

**Description**: This dataset consists of prompt and image URL pairs scraped from the LAION share-dalle-3 discord channel. The purpose is to collect image-text pairs for research and exploration.

# The Linked Open Data Cloud

## LINKED DATA
★ On the web, open license
★★ Machine-readable data
★★★ Non-proprietary format
★★★★ RDF standards
★★★★★ Linked RDF

**IS YOUR DATA 5 ★ ?**

Legend
- Cross Domain
- Geography
- Government
- Life Sciences
- Linguistics
- Media
- Publications
- Social Networking
- User Generated



https://lod-cloud.net/

# BIO2RDF

*Linked Data for the Life Sciences*

**Bio2RDF is an open source project that uses semantic web technologies to make it easier to reuse biomedical data.**

**It provides Linked Data and a queryable RDF knowledge graph.**

chemicals/drugs/formulations,
genomes/genes/proteins, domains
Interactions, complexes & pathways
animal models and phenotypes
Disease, genetic markers, treatments
Terminologies & publications

- **Since 2007, last updated 2014**
- **30+** biomedical data sources
- **10B+** interlinked statements
- NCBI, EBI, SIB, DBCLS, NCBO, and many others (chem2bio2rdf) produce this content

Belleau et al. JBI 2008. 41(5):706-716.
Callahan et al. ESWC 2013. 200-212

# the Triple as
# a base unit of knowledge representation

predicate

| subject | → | object |

"diclofenac is a drug"

is a

| diclofenac | → | drug |

# formalization

"diclofenac is a drug"



1. Use RDF
2. Assign/reuse identifiers
3. Use or develop vocabularies

**RDF N-Triples format (standardized, machine interpretable):**
<https://bio2rdf.org/drugbank:DB00586>
      <http://www.w3.org/1999/02/22-rdf-syntax-ns#type>
           <https://bio2rdf.org/drugbank_vocabulary:Drug> .

# Biomedical Linked Data

**About: Diclofenac [drugbank:DB00586]**  Goto  Sponge  NotDistinct  Permalink
An Entity of Type : http://bio2rdf.org/drugbank_vocabulary:Small-molecule, within Data Space : bio2rdf.org associated with source document(s)

Type: [ Small molecule [drugbank_vocabulary:Small-molecule] ⌄ ]  [ New Facet based on Instances of this Class ]

http(s) identifier

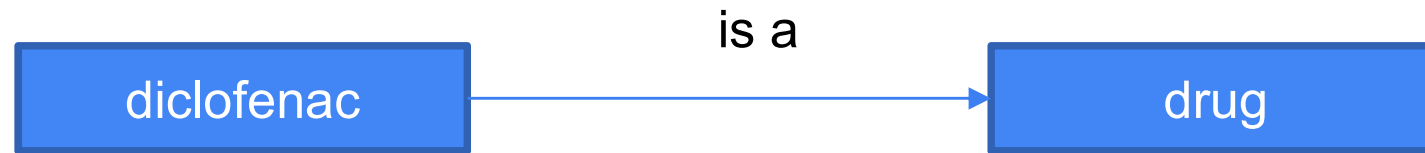| Attributes | Values |
|---|---|
| rdf:type | Drug [drugbank_vocabulary:Drug] |
| | drugbank resource [drugbank_vocabulary:Resource] |
| | Small molecule [drugbank_vocabulary:Small-molecule] |
| rdfs:label | Diclofenac [drugbank:DB00586] |
| rdfs:seeAlso | http://www.drugs.com/cdi/diclofenac-drops.html |
| | http://www.rxlist.com/cgi/generic/diclofen.htm |
| | http://www.drugbank.ca/drugs/DB00586 |
| sameAs | Diclofenac [drugbank:DB00586] |
| dcterms:title | Diclofenac |
| dcterms:description | A non-steroidal anti-inflammatory agent (NSAID) with antipyretic and analgesic actions. It is primarily available as the sodium salt. [PubChem] |
| dcterms:identifier | drugbank:DB00586 |
| void:inDataset | http://bio2rdf.org/kegg_resource:bio2rdf.dataset.kegg.R3 |
| | http://bio2rdf.org/drugbank_resource:bio2rdf.dataset.drugbank.R3 |
| | http://bio2rdf.org/pharmgkb_resource:bio2rdf.dataset.pharmgkb.R3 |
| Bio2RDF identifier | DB00586 |
| Bio2RDF namespace | drugbank |
| Bio2RDF uri | http://bio2rdf.org/drugbank:DB00586 |
| identifiers.org URI | Diclofenac [drugbank:DB00586] |
| x bindingdb [drugb...lary:x-bindingdb] | http://bio2rdf.org/bindingdb:13066 |
| x chemspider [drug...ary:x-chemspider] | http://bio2rdf.org/chemspider:2925 |
| x pdb [drugbank_vocabulary:x-pdb] | http://bio2rdf.org/pdb:DIF |
| x pubchemcompound ...-pubchemcompound] | http://bio2rdf.org/pubchem.compound:3033 |
| x pubchemsubstance...pubchemsubstance] | http://bio2rdf.org/pubchem.substance:46504644 |
| absorption [drugba...ulary:absorption] | absorption for drugbank:DB00586 [drugbank_resource:af3a8b347e732d3c3b48a5428a6160e0] |
| affected organism ...ffected-organism] | Humans and other mammals [drugbank_vocabulary:e1e572616d493b2affcc653e19cbcd21] |
| brand [drugbank_vocabulary:brand] | [drugbank_vocabulary:e78186eb12eeaebda8a530a67513beea] |
| | Aclonac [drugbank_vocabulary:6856a4532f20c29e5c1eac027b253a36] |
| calculated propert...lated-properties] | Traditional IUPAC Name: diclofenac from ChemAxon [drugbank_resource:calculated-properties-DB00586-10] |
| | Molecular Weight: 296.149 from ChemAxon [drugbank_resource:calculated-properties-DB00586-11] |
| | Monoisotopic Weight: 295.016684015 from ChemAxon [drugbank_resource:calculated-properties-DB00586-12] |
| | SMILES: OC(=O)CC1=CC=CC=C1NC1=C(Cl)C=CC=C1Cl from ChemAxon [drugbank_resource:calculated-properties-DB00586-13] |
| | Molecular Formula: C14H11Cl2NO2 from ChemAxon [drugbank_resource:calculated-properties-DB00586-14] |
| | »more» |

semantically typed

has detailed provenance

linked to other resources

rich descriptions

# Bio2RDF is a (Syntactically) Interoperable Biomedical Knowledge Graph

**Information Retrieval:** Phenotypes of knock-out mouse models for the targets of a selected drug



**Exploration**: drug-target-disease networks



https://doi.org/10.7717/peerj-cs.106

**Reproducible ML**: new uses for existing drugs

**Custom Knowledge Portal: EbolaKB**



https://doi.org/10.1093/database/bav049



https://doi.org/10.7717/peerj-cs.281

# Knowledge Collaboratory (for small data)

An *AI-powered* (NLP model or GPT) web user interface to annotate biomedical text (NER, RE), create standard-compliant statements (BioLink model) that can be made publicly available as author-signed nanopublications.

📝 **collaboratory.semanticscience.org/annotate**



Clonazepam `ChemicalEntity` is useful alone or as an adjunct in the treatment of the Lennox-Gastaut syndrome `DiseaseOrPhenotypicFeature` ( petit mal variant `DiseaseOrPhenotypicFeature` ), akinetic and myoclonic seizures `DiseaseOrPhenotypicFeature` .

2. Define the statements that represent the assertions made in the text, you can add properties to provide more context:

| Subject | Predicate | Object |
|---|---|---|
| CLONAZEPAM (PUBCHEM.COMPOUND:2802) | treats (biolink:treats) | lennox gastaut syndrome (MONDO:0016532) |

➕ Add a property to this statement

# BioLink Model

A data model to structure (qualified) biological associations.



"Bisphenol A results in decreased degradation of ESR1 protein" - A Statement where the effect has a direction (decreased)

```
{
"id": "e0",
"category": "biolink:ChemicalAffectsGeneAssociation",
"subject": "CHEBI:16811"  # Bisphenol A,
"predicate": "biolink:affects",
"qualified_predicate": "biolink:causes",
"object": "NCBIGene:2099"  # ESR1,
"object_aspect_qualifier": "degradation",
"object_direction_qualifier": "decreased"
}
```

https://biolink.github.io/biolink-model/

Technology to **publish assertions using RDF**

Contains RDF triples to specify the assertion, its provenance, and digital object metadata

Digitally signed by agent

TrustyURI hash to provide globally unique, persistent, immutable, verifiable identifier and payload

## Nanopub Monitor

32 **nanopub services** running on approximately 10 **distinct servers**

# Translator finds associations between drugs, genes, and diseases

## Select a question and enter a search term to get started

| What drugs may treat conditions related to ... ⌄ | ⚙ Log In to Enter a Search Term | → |

**Examples**
Choose a different question for more examples. Run a new search with these terms for the most up-to-date results.                    ⌄

| **Bethlem Myopathy** | **Breast Cancer** | **Chronic Obstructive Pulmonary Disease** | **Ehlers-Danlos Syndrome** |

## Translator Workflow



Diseases
Chemicals
Drugs

←

Favorite Results

←

View Evidence

**Select a relationship to explore**        **Review and select your favorite results**        **Analyze evidence in the workspace**

# FAIR

Robust, Reproducible, Explainable Predictions

# Neurosymbolic AI

NAI aims to combine symbolic reasoning methods (logic-based reasoning & rules) with sub-symbolic methods (neural networks, deep learning) to create models with high predictive performance and explanability.

Specifically:

- **Integrate knowledge** from different modalities
- Perform **complex reasoning** (e.g. deduction, induction, synthesis)
- **Learn from examples/small data and big data**
- Robust to noise and nonsense
- Handle cases **out of the learning distribution**
- Offer **explanations** (e.g. causal account of the phenomenon) and **justifications** (the evidence that supports the claims)

# Predict new drug applications in a documented and reproducible manner

Mol Syst Biol. 2011; 7: 496.
Published online 2011 Jun 7. doi: 10.1038/msb.2011.26

PMCID: PMC3159979

**PREDICT: a method for inferring novel drug indications with application to personalized medicine**

Assaf Gottlieb,[1] Gideon Y Stein,[2,3] Eytan Ruppin,[1,2] and Roded Sharan[a,1]

**AUC 0.90 across all therapeutic indications**

**Scripts not available. Feature tables available.**
**Not reproducible!**



**Towards FAIR protocols and workflows: the OpenPREDICT use case**

Remzi Celebi[1,*], Joao Rebelo Moreira[2,*], Ahmed A. Hassan[3], Sandeep Ayyar[4], Lars Ridder[5], Tobias Kuhn[2] and Michel Dumontier[1]

**Result: ROCAUC 0.83**

Celebi R, Rebelo Moreira J, Hassan AA, Ayyar S, Ridder L, Kuhn T, Dumontier M. 2020. Towards FAIR protocols and workflows: the OpenPREDICT use case. PeerJ Computer Science 6:e281 https://doi.org/10.7717/peerj-cs.281

36

# Explainable AI

$$Score(d_r, d_i) = \max_{d_r', d_i' \neq d_r, d_i} \sqrt{S(d_r, d_r') \times S(d_i, d_i')} \qquad (2)$$

- XAI methods such as SHAP provide information about feature importance for the model and in individual predictions
- When applied to OpenPredict, **it's too complicated understand** the contributions of derived features
- However, it is clearer when using a single feature predictor



**Fig. 7**: Global explanations of OpenPREDICT model as mean absolute of SHAP value.



**Fig. 9**: Global explanations of XPREDICT Single model as mean absolute of SHAP values.

# Graph Representation Learning

We want to automatically discover effective representations needed for classification from data.

In graph representation learning, we encode the topology, node attributes, and edge information into low-dimensional vectors (or embeddings).

These vector can then be used as features to train classifiers for link prediction, node classification, graph classification, etc

# Graph Neural Networks

Graph Neural Networks (GNNs) iteratively update node representations by aggregating features from neighbouring nodes and possibly edges.

Several methods (e.g. Saliency Maps) exist to extract a model-wise explanation for link prediction, node/graph classification.



Fig. 2: Interconnections among entities in DRKG [19].

GNN

+

XAI



Fig. 3: Explanatory subgraph consisting of only drugs and diseases for (Memantine, treats, Alzheimer) using Graph Attention Network and Saliency Maps.

**But I don't find these explanations salient at all.**

**They lack the sophistication of a reasoned explanation for the predicted phenomenon**

# Building better explanatory subgraphs by combining GNNs + XAI + KGs

Use **Graph Neural Networks** to capture **semantics**, **graph structure** and **relationships between nodes**

Apply **Saliency Maps** on predictions made by GNNs to **identify relevant nodes** for a specific prediction; this provides valuable insights into the graph's topology and highlights the most important components

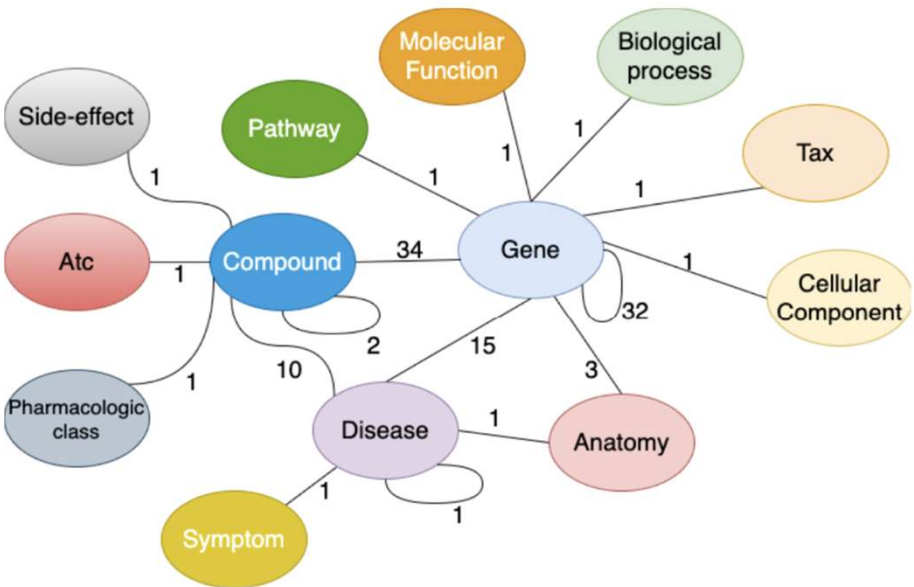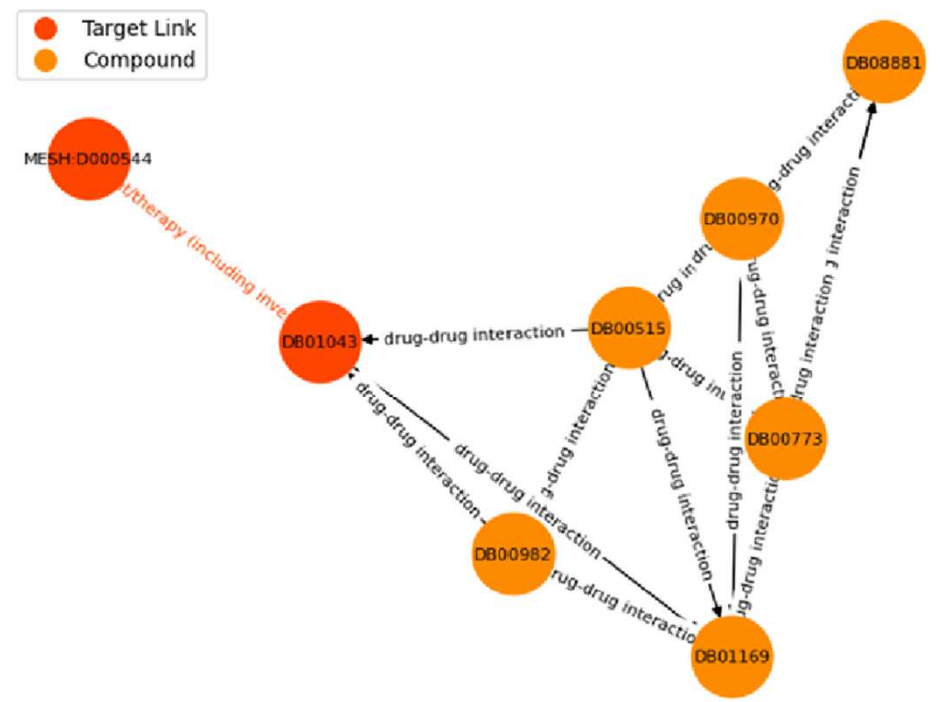Saliency Maps assign a score to each node in the network, which can be used to **rank paths** involving genes, pathways, diseases, and compounds



---

**Algorithm 1** Algorithm used to generate explanations.

1: **function** GENERATEEXPLANATORYSUBGRAPH($SM\_scores, k$)  ▷ where $SM\_scores$ - scores derived from Saliency Maps, $k$ - number of triples included in the explanation
2:    Let $g_1, g_2, ..., g_n$ be ranked gene entities based on $SM\_scores$
3:    RankedTriples = []
4:    **for** $g_i$ = 1 to $n$ **do**
5:        PathwayRel = ExtractRelations($g_i$, "ParticipatesIn", "Pathway")
6:        DiseaseRel = ExtractRelations("Disease", "Associates", $g_i$)
7:        CompoundRel = ExtractRelations("Compound", "BindsTo", $g_i$)
8:        RankedPathwayRel = RankRelations(PathwayRel)
9:        RankedDiseaseRel = RankRelations(DiseaseRel)
10:        RankedCompoundRel = RankRelations(CompoundRel)
11:        RankedTriples.append(RankedPathwayRel[:k], RankedDiseaseRel[:k], RankedCompoundRel[:k])
12:    **end for**
13:    ExplanatorySubgraph = BuildExplanatorySubgraph(RankedTriples)
        **return** ExplanatorySubgraph
14: **end function**

# Explanations for Alzheimer's Disease Treatments

By blocking the NMDA receptor, Memantine protects neurons from excitotoxicity caused by glutamate overstimulation, explaining how it reduces cognitive decline and improves overall function in Alzheimer's patients.

The primary goal of Alzheimer's drugs, including donepezil, is to maintain elevated acetylcholine (ACh) levels, thereby compensating for the loss of functioning cholinergic brain cells. Donepezil binds to acetylcholinesterase (AChE) and butyryl-cholinesterase (BChE), which are involved in ACh metabolism and thus are important for the cholinergic function in the brain



Figure 2: Explanatory subgraph for (Donepezil, treats, Alzheimer) including relationships from the original KG

# Compelling scientific justifications for ML model predictions

The goal is to provide *explanations* for predictions based on the scientific literature to complement model explanations.

Compelling justifications are a set of assertions derived from primary sources which provide evidence, both for and against, the stated assertion, and that reflects on the likelihood of any of those statements to be true and well supported.

We explore the feasibility of using RAG to create these compelling justifications from scientific text and structured knowledge graphs.

Model

KG

Prediction

PubMed®

XAI

RAG

LLM

Justification

Maastricht University

Stated fact = "Insulin human DOES NOT HELP to control hyperglycemia in diabetes mellitus."

#13

**Mechanisms of action:** Insulin is a hormone that regulates blood sugar levels by allowing cells to take up glucose from the bloodstream and use it for energy or storage. In people with diabetes, either the body does not produce enough insulin (type 1 diabetes) or the body's cells become resistant to its effects (type 2 diabetes), resulting in elevated blood sugar levels. Insulin delivery modalities, such as continuous subcutaneous insulin injection (CSII) and multiple daily injections (MDI), aim to replace or supplement the insulin that the body is not producing or using effectively.

**Evidence-based medicine:** The AHRQ systematic review of randomized clinical trials found that both CSII and MDI demonstrated similar effectiveness on glycemic control and severe hypoglycemia in children and adolescents with T1DM and in adults with T2DM, regardless of insulin delivery modality. In adults with T1DM, the review revealed a stronger reduction in hemoglobin A1c with CSII compared to MDI with low strength of evidence. However, one study heavily influenced these results. The review also highlights advances in insulin delivery and glucose-monitoring systems, practical methods to achieve tight glycemic control and strategies to minimize associated risks, as well as gaps in research and areas that need to be addressed in the future.

**Comparisons with other treatments:** In comparison to other treatments for diabetes, such as oral medications or bariatric surgery, insulin therapy has been shown to be effective in controlling blood sugar levels and reducing complications associated with high blood sugar, such as kidney damage, nerve damage, and blindness. However, insulin therapy can have side effects, such as hypoglycemia (low blood sugar), weight gain, and injection site reactions, which need to be carefully managed by healthcare providers.

**Fact:** Insulin human DOES help to control hyperglycemia in diabetes mellitus in adults and children with type 1 or type 2 diabetes when administered correctly using appropriate insulin delivery modalities.

**Maastricht University**

**Expert Evaluation:** These selected justifications were then compiled into a survey, and subsequently sent to a board-certified internal medicine physician and cardiac critical care specialist. The physician evaluated each justification on a 5-point Likert scale ranging from excellent to poor for the following metrics:

- **Factuality:** The degree to which the generated justification aligns with established medical knowledge.
- **Completeness:** The extent to which the generated text provides a comprehensive and accurate representation of the question posed, including the inclusion of contraindications as necessary.
- **Relevance:** The degree to which the generated justification directly addresses the specific query or clinical context, avoiding unnecessary information.
- **Consistency:** Ensuring that the generated justification is internally consistent and doesn't contradict itself or established medical knowledge. Accompanying these scores the physician also provided a brief description of the rationale behind their scores for each metric.

| |
|---|
| Lack of explicit citation |
| Plausible sounding hallucinations or inconsistencies |
| Failure to understand the prompt |
| Influenced generation due to retrieved context |
| Irrelevant retrieved context |
| Missing evidence |
| Missing minor details |

explain how memantine could be used to treat alzheimer's disease

Thought for 4 seconds ⌄

## Introduction

Alzheimer's disease is a progressive neurodegenerative disorder characterized by memory loss, cognitive decline, and behavioral changes. One of the pathological features of Alzheimer's is the overactivation of glutamate neurotransmission, leading to neuronal damage through excitotoxicity. Memantine is a medication used to manage moderate to severe Alzheimer's disease by targeting this glutamatergic system.

## Mechanism of Action

1. **NMDA Receptor Antagonism**: Memantine is an uncompetitive antagonist of the N-methyl-D-aspartate (NMDA) receptor, a subtype of glutamate receptor in the brain. Under normal conditions, glutamate binds to NMDA receptors to facilitate learning and memory. However, excessive glutamate release can lead to overactivation of NMDA receptors, causing increased calcium influx and neuronal damage.

2. **Regulation of Glutamate Activity**: Memantine selectively blocks NMDA receptors only during excessive stimulation. It attaches to the receptor's ion channel, preventing prolonged calcium influx without disrupting normal synaptic activity. This helps protect neurons from excitotoxicity while preserving essential neurotransmission.

**How should the quality of explanations for predictions be properly evaluated?**

**Trust and Accountability in Knowledge Graph-Based AI for Self Determination**
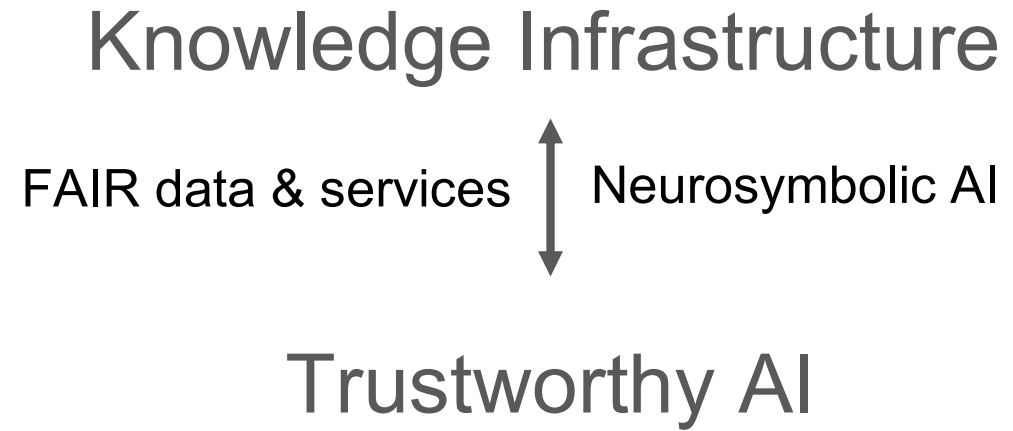**Dagstuhl Seminar 25051**

# Summary

The FAIR principles offer a high level set of requirements to increase productivity data-related tasks. *It is focused on making data available for software/machine agents*.

AI depends on ample amounts of high quality data – F**AI**R data – to train AI models. But *getting the right answer isn't enough* – we need (domain-based) explanations and compelling justifications to judiciously utilise this information.

**Neurosymbolic AI** combines the unreasonable effectiveness of machine learning predictions with the justificatory power of logic-based reasoning.

**LLMs** are rapidly evolving in their natural language understanding and human-like natural language generation – when combined with external knowledge sources (like KGs and NAI), hybrid systems show advanced reasoning capabilities.

# Towards Biomedical Neurosymbolic AI: From Semantic Knowledge Infrastructure to Explainable Predictions

Knowledge Infrastructure

FAIR data & services | Neurosymbolic AI

Trustworthy AI

**Michel Dumontier, PhD**

Distinguished Professor of Data Science
Founder and Director, Institute of Data Science
Department of Advanced Computing Sciences
Maastricht University

michel.dumontier@maastrichtuniversity.nl