

ANN: A platform to annotate text with Wikidata IDs

This manuscript ([permalink](#)) was automatically generated from [lubianat/ann_sprint@db79ceb](#) on September 10, 2020.

Authors

- **Tiago Lubiana**

 [0000-0003-2473-2313](#) ·  [lubianat](#) ·  [johndoe](#)

Computational Systems Biology Laboratory, University of São Paulo

Abstract

Individual contributions

All contributors participated in the discussions about the platform design, the directions of the project and on the draft of this manuscript.

Additionally:

- I.H.G.P , A.D.R, D.F, T.L, G.N.V and D.A. contributed to the manual annotation of scientific articles.
- S.W. and A.G worked on the analysis of Natural Language Processing software and on the integration of sci spaCy to Wikidata.
- D.A. made connections with Europe PMC and designed the frond-end mockup.
- A.D.R. organized the survey of the researchers' views
- G.N.V , D.F and A.D.R. discussed and wrote the community aspects on benefits of annotation and research incentives.
- B.A.P designed the logo.
- T.L, D.W, I.H.G.P worked on extracting patterns from manual annotations
- T.L proposed the initial idea, and oversaw the project.

Introduction

Most scientific findings are communicated via scientific papers. In these papers knowledge is encoded in language, which makes it difficult for computers to search for facts. Additionally, millions of articles are produced every year, which makes it virtually impossible to keep track of current biomedical knowledge without the aid of computational tools.

Annotation - linking words in scientific texts with external identifiers - is a ground step to automatically process biomedical knowledge. The annotation of biomedical articles is part of the work of biocurators, professionals dedicated to parse and make knowledge available on databases. These annotations, usually based on ontologies (sets of organized concepts) already power core platforms used by the research community, such as ZFIN and UniProt. (Note that the meaning of "annotation" used here is different from sharing written notes about ideas on the text such as proposed by [hypothes.is](#)).

Even though biomedical databases are extremely valuable, they are limited to specific subsets of human knowledge. It would be a colossal challenge to annotate text if we needed to look for the right ontology for each kind of concept.

Wikidata is a possible solution for this challenge. It is knowledge base that contains more than 80 million varied concepts: "p53", "malaria", "Douglas Adams", "Brazil" and much more. Moreover, anyone

can contribute with new concepts (and relations between them) to Wikidata, which makes it flexible enough to accommodate the vast amount of concepts used in research articles.

During the 2-day hackathon eLife Sprint 2020, we hot-started a project for annotating concepts in scientific articles to Wikidata, envisioning integration with Europe PMC's Annotation API. We brainstormed both technical and practical aspects of developing a tool to gather crowd-annotations of scientific concepts. Inspired by other scientific games (like Mark2cure, eterna and foldit), we designed a gamified interface for crowdsourcing scientific annotations. Additionally, we studied Natural Language Processing approaches for extracting scientific entities, and assembled a series of perspectives on how to implement such an annotation tool in the current research environment.

This document contains reports on the different branches of the project, coupled to thoughts of the participants on how to achieve the overarching goal of annotating all scientific text. Given the short time for the event, some parts of the report are not completely structured. Nevertheless, we believe that they can be useful for accessing the development of the project.

Tasks

Tasks developed at the hackathon At the hackathon we worked on the following tasks:

- Deploy a pilot survey of the researchers' views on annotating of scientific texts
- Develop code for selecting candidate concepts for annotations (using the NLP package scispacy) Review works related to annotation of texts and how they relate to the project.
- Analyse the incentive structure and how to engage researchers in the annotation of scientific texts
- Manually annotate a sample of biomedical publications to Wikidata IDs and compare them to currently annotated pieces of text on Europe PMC
- Extract patterns related to the availability of annotated concepts in Wikidata.
- Design a logo and a name for the project
- Design an user interface and prepare a mockup "front-end".

Why is annotation important?

Two main contributions of annotations are to clarify the meaning of texts and enable programmatically processing.

Clarify scientific prose

- **Disambiguation of concepts and abbreviations:** Words can have different meanings in different contexts. With annotations, selected words in a scientific article or abstract are linked to a concept in a knowledge base (e.g. Wikidata) and can be clearly disambiguated. Ultimately, this helps readers better understand the content of scientific articles.
- **Science communication:** through an enriched layer, annotations help convey complex information to the general public. It is possible to implement a "hover" function (similar to that in Wikipedia), where the user could hover above a word and you see the definition of a concept. Programmatically process articles Semantic enrichment: annotations enrich a scientific article with an additional layer of machine-readable information, providing more in-depth information about a concept or linking it to other sources of information.

The semantic enrichment layer unleashes the biomedical knowledge constrained in biomedical articles to the world of the semantic web. The connections of concepts, then, can impact a number of different aspects of research:

Programmatically process articles

- **Semantic enrichment:** annotations enrich a scientific article with an additional layer of machine-readable information, providing more in-depth information about a concept or linking it to other sources of information.

The semantic enrichment layer unleashes the biomedical knowledge constrained in biomedical articles to the world of the semantic web. The connections of concepts, then, can impact a number of different aspects of research:

- **Integration of different sources of information:** annotations can help us find information related to any given concept, regardless of its source. From the perspective of the researcher, this can improve the visibility of their work, making it reachable for the ones interested in the area.
- **Improvement of document classification:** annotation can help automated document classification, making it easier to search these documents.
- **Search for complex questions:** Text annotated with Wikidata IDs (semantically enriched) are readable by computers. This enables better discovery mechanisms (and not just left as words). Moreover, Wikidata is a knowledge graph, and concepts are linked to each other. That makes it possible to leverage the collective knowledge embedded in the graph to make powerful queries, such as: "Which articles produced by alumni of my university mention drugs that block NMDA receptors?" , "Which cell lines are used for research that deals with respiratory viruses"?
- **Improve openness of research:** annotations in a paper that link to an open knowledge base increase the openness of an article in accordance to the FAIR principles of Findability, Accessibility, Interoperability and Reusability. Our work focus on annotations compatible with the EuropePMC API (which uses the W3C standard and encodes annotations in a RDF-compatible format), therefore making annotations quickly available via the API itself and wrappers, such as the R package europe PMC.

Related softwares for annotation

The task of connecting mentions in scientific texts to identifiers in databases has been researched in both the biomedical and the natural language processing communities. There have been many approaches to automate this task. This is different from our human-in-the-loop approach. Nevertheless different components can help us to select candidates that we present to the annotators.

Of note, there have also been previous approaches for annotating documents using Wikipedia [\[1\]](#)

Overview

- Open Tapioka
 - Recognizes entities and links them to Wikidata, but only for person, location organization
 - Could be retrained with a subset of Wikidata that would contain only biomedical entities
- Sci Spacy
 - Entity detection and linking to UMLS ID
 - Linked WikiText-2 Project
 - Does Entity identification, Annotation with Wikidata entries.
 - The project utilizes Stanford CoreNLP
- Doccano
 - Does Entity identification, Sentiment analysis.
- BERN

- Uses contextualized word embeddings, which might have higher accuracy than sci Spacy
- Does Entity Identification, Entity typing

Open Tapioka

- [github repo](#)
- [online demo](#)
- [paper](#)
- [documentation](#)

System description

Input is a sentence. Output is the sentence, with all persons, locations and organizations linked to their respective Wikidata identifier. The system is trained solely on Wikidata. The authors use occurrence statistics of concepts in Wikidata and in text to compute the likelihood of a certain word in the text linking to a certain Wikidata item (e.g. "Barack Obama" linking to "Q76"). To take context into account independently computed local features are propagated along a Markov chain. The authors claim that this system is lightweight and easy to retrain, and therefore easily adapts to the frequent changes of Wikidata. They say that restricting their system to only people, organizations and locations enabled them to do well without using any other data but Wikidata, while other approaches do rely on additional text from Wikipedia.

How feasible is this system for our project

Using only Wikidata to train the system is a good asset, because this might keep training times low. The authors are right in claiming that their system is lightweight: It does not use word embeddings or extensive language models but derives the necessary information about word similarities from Wikidata itself. The author states that this approach worked for relatively common entities, so it is unclear whether we can adapt it to less common biological entities. Testing the system on the cited website gave reasonably good results for less common names of people and cities, but was prone to misinterpret words that were not people, organization or locations as such (e.g. in the sentence "Banks are often closed." the words "Banks" and "closed" were linked to locations). The documentation seems generally very good.

Questions/Answers:

- How does the system pick out only people, locations and organizations? This is done before training, by using only entities of those categories in the training data set. (see documentation [here](#))
- How easily can this be changed in the code? If we had a dump of Wikidata containing only biological entities we could use it to train on as described, no changes to the code itself needed!

Sci Spacy

- [github link](#)
- [online demo](#)
- [paper](#)

System description

The input is a sentence. The output is a sentence with the biomedical entities in that sentence annotated with canonical names, concept IDs and TUI(s).

How feasible is this system for our project:

This looks nearly perfect for candidate generation.

Questions:

- Are the IDs provided there in any way meaningful for linking them to Wikidata? Yes! (See section below)

Doccano

- [code base](#)
- [demo](#)

Comments: - Entity identification; Sentiment analysis.

Linked WikiText-2 Project

- [codebase](#)
- [demo](#)
- [backbone](#)

Entity identification; Annotation with Wikidata entries; the project utilizes Stanford CoreNLP

BERN

- [code base](#)
- [demo](#)

Uses contextualized word embeddings, which might have higher accuracy than sci Spacy Does Entity Identification, Entity typing

Applicability of the sciSpacy tool

The input for the software backend are abstracts of scientific articles that are loaded from Europe PMC using the Europe PMC API. We then use sci Spacy to detect entities in the abstract. Sci Spacy annotates those entities with their ID in the Unified Medical Language System (UMLS).


Notably, 26 thousand items in Wikidata have an UMLS ID, which allows to link the items that were detected by sciSpacy to be connected to Wikidata. We pinpoint a couple challenges:

The pre-trained scispacy models are unable to identify the entity if it has conjunctions and prepositions in it. To improve the entity detection performance the model needs to be retrained using manually curated scientific word lists. Sci spacy does not use contextualized word embeddings, which impacts the precision of retrieved entities (the model employed in sci Spacy is derived from this [reference\[2\]](#)).

Other approaches use contextualized word embeddings for detecting and normalizing biomedical entities (such as <https://bern.korea.ac.kr/>). They could be used in the project in addition to scispacy.

Besides UMLS IDs, sciSpacy also can match concepts a number of MeSH IDs, which can also be linked to Wikidata items.

We wrote code in a Google Colab notebook to concatenate the Europe PMC API with sciSpacy and Wikidata. After retrieving the abstract of an article via its PMID, the function extracts relevant concepts via sciSpacy and match the ones with PMIDs to Wikidata. The output of the pipeline is depicted in the Figure [??] and the code is available in the [project github repository](https://github.com/lubianat/ann)

 `get_pmcid_annotations("32543932")`

	label	start_pos	end_pos	mesh_id	meshid_match_score	arbitrary_wdata_id	wdata_id
0	Treatment	0	9	D013812	0.856167	Q179661	Q179661
1	hypomethylating agents	15	37	None	NA	None	None
2	HMA	39	43	None	NA	Q838699	None
3	azacitidine	45	56	None	NA	Q416451	None
4	decitabine	60	70	D000077209	0.709714	Q1181878	Q1181878
5	standard of care	86	102	D059039	0.702054	None	Q7598360
6	high risk myelodysplastic syndromes	107	142	None	NA	None	None
7	MDS	144	148	D009190	0.751488	Q57403622	Q954625
8	associated with	157	172	None	NA	None	None
9	low rates	173	182	None	NA	None	None

Figure 1: Concatenation of the Europe PMC API to Wikidata and sciSpacy

Software Frontend

A mockup of the frontend is available in Figure 1 and the rules used for the user interface design are shown in Figure 2. The main idea is to make annotations made with ANN fun for annotators. Users will be able to search for a publication by PMID or title. ANN boxes will be filled in with title and abstract. Annotations will work in a task manner where annotators will be asked to annotate a type of entity or sentence and be rewarded with ANN badge points.

When selecting a text for annotation, a window will popup and users will be able to select terms from Wikidata. The annotations will be saved in a format that is compatible with Europe PMC annotations submission system, which would add Wikidata annotations to the Europe PMC SciLite annotations features. Users would be able to login with an ORCID account and ANN would allow them to claim their annotations work to their ORCID account.



Ann Smith 
You have earned
50 

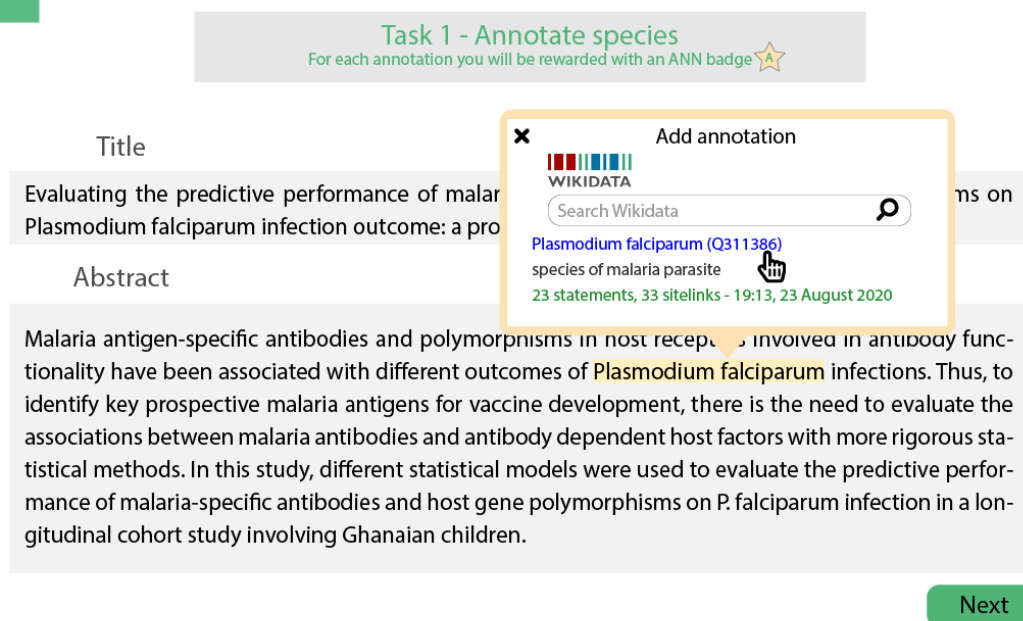


Figure 1: Mockup prototype of the frontend of ANN



Ann Smith 
You have earned
50 

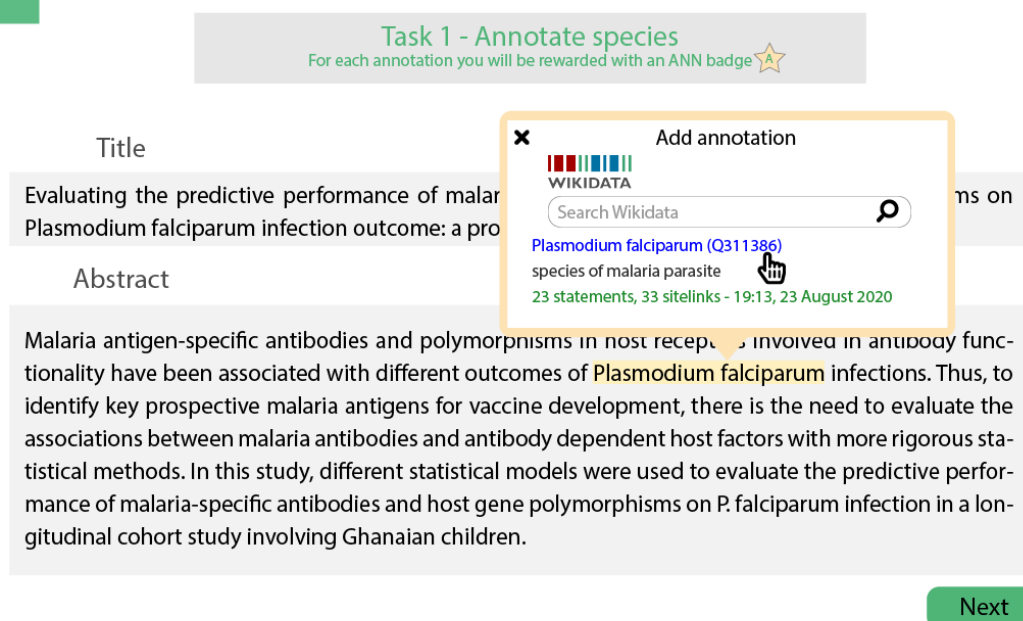


Figure 2: [Miro](#) board for brainstorming of the functionalities of the user interface.

Manual annotations: case report

We manually annotated a sample of biomedical publications (some of which from eLife) to extract information from their abstracts as a prototype for the platform. These annotations were made mostly by people from a research-lab background, which were presented to Wikidata during this project. In that sense, they mimic what a researcher could annotate when using the ANN platform.

The annotations and notes are available in the attached table: [Manual Annotations Spreadsheet](#).

The following observations were made:

Annotation of PMID 19268344 [3]

“Ankyrin G” was not found on Wikidata, but only “ankyrin”. The [ankirin](#) item on Wikidata refers to a [protein family](#), and not to the protein itself. Later, we noticed that “Ankyrin G” was actually present on Wikidata ([here](#) or [here](#)), but at the moment the exact term “Ankyrin G” was not listed as one alias.

- A1 - Concepts might be present on Wikidata, but not with the exact wording they appear in scientific texts.
- A2 - The abstract of the text does not mention if the study is dealing with humans or mouse cells! Digging into the text, it mentions that “Dissociated hippocampal neurons were prepared from embryonic E18 rats.” The actual annotation would be [this](#), but this was impossible to tell given the abstract! Annotation would solve this, and the authors would be the best to disambiguate.

Given A3, There are some things that should always be annotated, such as the species. Which kind of things should be always annotated? How should they be required? Things such as the species of the organisms used in the study, and the techniques used to assess the question. Maybe if journals require a semantic abstract, a structured abstract with semantically enriched knowledge, some of this problems would be bypassed.

The concept of “[hippocampal neurons](#)” was not identified separately on Wikidata, just the concepts for “[hippocampus](#)” and “[neuron](#)”. Notably, even these two concepts are represented in an species-independent way. The ideal annotations would be precise for the taxon of interest.

- A3 - Compound concepts might be split on Wikidata on its building parts. However, as per the view of the annotator, “hippocampal neuron” can be understood as an individual concept. This illustrates an open challenge of annotation: do we annotate the parts, or do we annotate the complete concept?

AIS, the acronym for “Axon Initial Segment” did not have a string-match on Wikidata. Notably, the concept is [present on Wikidata](#), but not the wording of the acronym, making this an instance of the problem A1 (lack of exact name).

- A4 - The AIS acronym is present on Wikidata but for many different things (like the Australian Institute of Sport and in situ pulmonary adenocarcinoma). Acronyms are tricky, and might need disambiguation before matching to the database. Some NLP programs like spacy have acronym disambiguation modules. Adding a concept to Wikidata in the perspective of a new user

Adding a concept to Wikidata in the perspective of a new user

The concept of “[hippocampal neurons](#)” was created on Wikidata by a new editor. However, the creation of an item in an ontology is a challenging task. The terms used are not common for the biological research workflow. These are the perceptions of one of the team members, when creating an Wikidata item for the first time:

- LABEL: main annotation for the concept. This is what would be used from SpaCY,
- DESCRIPTION: what the annotation refers to. As simple as possible. Simple words, it can reference other concepts.
- ALSO KNOWN AS/ alias : other concepts that mean the same thing (I wonder whether this could be “counted” as one, if referencing throughout an abstract). For example, when writing an abstract you use “different words” for the same thing to avoid repetition, would this be understood as one same concept if they are linked ? (in multiple languages): All this can be done in multiple languages:
- When adding statements:
 - you can specify if the concept is a subcategory of a different concept: SUBCLASS OF
 - you can specify if the concept is an example of a different concept: INSTANCE OF
 - You can also specify if the concept is PART OF: I take this to be as part of a bigger system, not necessarily immediate subclass. But this can be tricky to distinguish I think.
 - You can add an IMAGE to better describe it." _

This is the first report (as of our knowledge) of a biology researcher perspective when faced with the task of adding a new item to Wikidata.

Given the complexity of ontological modelling, the report shows that the task is feasible. Adding a concept without training in semantic technology makes Wikidata a powerful tool, as the barrier of contribution is much lower than the one for current ontologies. The entries can be, then, adjusted later by the community, if necessary.

Annotation of PMID 31254741 [4]

This article is related to drosophila research. Drosophila genes have frequently funny names, which might be mistaken for other entities. For the gene [frizzled](#), the system worked nicely. We could find the Wikidata entry by typing “frizzled” but the official name is “fz”. The [protein](#) is also present, and the researchers would have to choose if they are talking about the gene or its product.

The Van Gogh protein, on the other hand, does not show up on Wikidata, just the [dutch post-impressionist painter](#). For this case, the protein is not on Wikidata, but the [protein-coding gene](#) is, but it gets “buried” amidst the references to the painter.

- B1 - Many entities may have the same name, even when dealing with full words (and not acronyms). This ambiguity cannot be solved by looking at the word, and it might be solvable (at least partially) by looking at the context.
- B2 - A biologically similar entity might be present, but in a slightly different way (for example, gene entry when talking about a protein). This is a more general case of problem A2, of disambiguation of gene names for different species.

Annotation of PMID 31909712 [5]

For this article, PCP (planar cell polarity) was not found on Wikidata, just other references for the same acronym (instance of case A4). In this case, even searching for the full concept does not lead to a hit.

- C1- The concept of interest may be completely missing.

Of note, many scientific articles catalogued on Wikidata mention “planar cell polarity” in their titles, which might be confusing for annotation. This is a good indication that it would be useful to create the concept on Wikidata, at the very least to link to the article items via the [main subject property](#).

- C2 - The concept itself is not present, but there are scientific articles with the concept in the title.

For [N-cadherin](#), the label on wikidata was “cadherin 2”, but “N-cadherin was listed as an alias. This is not a problem, but it seems to be a possible source of confusion, so it is worth mentioning it.

- C3 - The main label on Wikidata is not the one used on the article, but the name used by the article is present as an alias.

Conclusion of manual annotations

The case study of manual annotations is useful to find patterns of problems in database matching. Notably, we had time only to analyse 3 of the many manual annotations made during the eLifeSprint 2020. In that way, the spreadsheet of manual annotations represent a rich resource for further exploration of the details related to annotation of biological concepts to Wikidata.

References

1. Annotating Documents by Wikipedia Concepts

Peter Schönhofen

Institute of Electrical and Electronics Engineers (IEEE) (2008-12) <https://doi.org/bqfkzk>

DOI: [10.1109/wiat.2008.56](https://doi.org/10.1109/wiat.2008.56)

2. Neural Architectures for Named Entity Recognition

Guillaume Lample, Miguel Ballesteros, Sandeep Subramanian, Kazuya Kawakami, Chris Dyer

Association for Computational Linguistics (ACL) (2016) <https://doi.org/gf2k56>

DOI: [10.18653/v1/n16-1030](https://doi.org/10.18653/v1/n16-1030)

3. A selective filter for cytoplasmic transport at the axon initial segment.

Ai-Hong Song, Dong Wang, Gang Chen, Yuju Li, Jianhong Luo, Shumin Duan, Mu-Ming Poo

Cell (2009-03-05) <https://www.ncbi.nlm.nih.gov/pubmed/19268344>

DOI: [10.1016/j.cell.2009.01.016](https://doi.org/10.1016/j.cell.2009.01.016) · PMID: [19268344](https://pubmed.ncbi.nlm.nih.gov/19268344/)

4. Experimental and Theoretical Evidence for Bidirectional Signaling via Core Planar Polarity Protein Complexes in *Drosophila*.

Katherine H Fisher, David Strutt, Alexander G Fletcher

iScience (2019-06-18) <https://www.ncbi.nlm.nih.gov/pubmed/31254741>

DOI: [10.1016/j.isci.2019.06.021](https://doi.org/10.1016/j.isci.2019.06.021) · PMID: [31254741](https://pubmed.ncbi.nlm.nih.gov/31254741/) · PMCID: [PMC6610702](https://pubmed.ncbi.nlm.nih.gov/PMC6610702/)

5. Vangl2 acts at the interface between actin and N-cadherin to modulate mammalian neuronal outgrowth.

Steve Dos-Santos Carvalho, Maite M Moreau, Yeri Esther Hien, Mikael Garcia, Nathalie Aubailly,

Deborah J Henderson, Vincent Studer, Nathalie Sans, Olivier Thoumine, Mireille Montcouquiol

eLife (2020-01-07) <https://www.ncbi.nlm.nih.gov/pubmed/31909712>

DOI: [10.7554/elife.51822](https://doi.org/10.7554/elife.51822) · PMID: [31909712](https://pubmed.ncbi.nlm.nih.gov/31909712/) · PMCID: [PMC6946565](https://pubmed.ncbi.nlm.nih.gov/PMC6946565/)