

PKG 2.0 Dataset Description

The PubMed Knowledge Graph is available at <https://pubmedkg.github.io> and is updated annually. The original document of each field description of PubMed can be found at: https://www.nlm.nih.gov/bsd/licensee/elements_descriptions.html.

PKG can be imported into relational databases such as MySQL, making it easier to uncover patterns and relationships within the PKG dataset.

For instance, to query all papers published by a specific author, one can use the following SQL code:

```
SELECT PMID FROM C02_Link_Papers_Authors WHERE AID=1929791 # Assume the author's ID is 1929791
```

Query all patents that cite a specific paper:

```
SELECT * FROM `C16_Link_Patents_Papers` WHERE PMID=9923457 # Assume the paper's ID is 9923457.
```

Query all biomedical entities of a specific clinical trial:

```
SELECT * FROM `C13_Link_ClinicalTrials_BioEntities` WHERE nct_id='NCT00000102' # Assume the paper's ID is NCT00000102.
```

Query the bibliometric information of a specific author:

```
SELECT * FROM `C07_Authors` WHERE AID=666666 # Assume the author's ID is 666666.
```

Query all papers funded by a specific fund:

```
SELECT * FROM `B03_Link_Papers_Projects` WHERE PROJECT_NUMBER='Z01HD000364' # Assume the project's ID is Z01HD000364.
```

The tables starting with the beginning of A are the original data tables of PubMed. The tables with the beginning of B are the tables associated with external data sources to provide a support for PKG2.0. The tables with the beginning of C are synthesized by extracting key information from the tables A and B for subsequent statistics and calculations.

1. A01_Articles

Specific information for each article.

Table 1 Data Description for records of A01_Articles

Column Name	Description
id	Automatically incremented identifier
PMID	Unique ID assigned by PubMed to identify articles.
PMID_Version	PMID version number, which is the record corresponding to all PMID versions in this table. That is, if the PMID version number is different, there are different records in this table, and the id field of each record is different. In other related article tables (such as C01_Article_simple), only the record with the highest version number is kept.
MedlineCitation_Owner	Organization responsible for creating and verifying citations, including: NLM, NASA, PIP, KIE, HSR, HMD, SIS, NOTNLM
MedlineCitation_Status	It is the stage of the article. There are seven possible values: Completed In-Process PubMed-not-MEDLINE In-Data-Review Publisher MEDLINE OLDMEDLINE
Journal_JournalIssue_PubDate_Year	Publication year of the current article.
Journal_JournalIssue_PubDate_MedlineDate	Publication year of the current article. If <i>Journal_JournalIssue_PubDate_Year</i> is null, the publication year of the current article is recorded in this field

2. A02_AuthorList

Specific information for each author.

Table 2 Data Description for records of A02_AuthorList

Column Name	Description
PMID	Unique ID assigned by PubMed to identify articles.
Au_Order	Author order of the current authors in the author list of current articles.
AuthorNum	Co-author number of the current article
Vetle_aid	Unique author ID allocated by Vetle.
S2ID	Unique author ID allocated by Semantic Scholar
AID	Unique author ID (The final author disambiguation result. S2id is the main source, supplemented by other disambiguation results)

3. A03_KeywordList

Article keyword information: keyword information in this table is provided by the data producer.

Table 3 Data Description for records of A03_KeywordList

Column Name	Description
PMID	Unique ID assigned by PubMed to identify articles.
Keyword	Keywords of the current article

4. A04_Abtract

The abstract of each article.

Table 4 Data Description for records of A04_Abtract

Column Name	Description
PMID	Unique ID assigned by PubMed to identify articles.
AbstractText	The abstract of the current Article

5. A05_GrantList

Grants details of each article.

Table 5 Data Description for records of A05_GrantList

Column Name	Description
PMID	Unique ID assigned by PubMed to identify articles.
GrantID	Funding identifiers, including research grant numbers or contract numbers (or both) that are financially supported by the US Public Health Service or any agency of the National Institutes of Health (NIH)
Project_Number	The NIH project number corresponding to the funded article. Each article may be funded by multiple projects. Multiple project numbers are separated by ",".

6. A06_MeshheadingList

Mesh Heading details of each article. Mesh Heading refers to the NLM control vocabulary and medical subject heading (MeSH®), which is used to characterize the content of the articles represented by MEDLINE citations.

Table 6 Data Description for records of A06_MeshheadingList

Column Name	Description
PMID	Unique ID assigned by PubMed to identify articles.
DescriptorName	Descriptors extracted from each article
DescriptorName_MajorTopicYN	If the MeSH descriptors assigned to the article is the Key word of the article, the value is Y, otherwise it is N.
DescriptorName_UI	Identify MeSH's unique encoding for each descriptor and qualifier
QualifierName	Qualifiers, including numbers and words

7. A07_SupplMeshList

Supplementary conceptual terms and protocol terms for each article.

Table 7 Data Description for records of A07_SupplMeshList

Column Name	Description
PMID	Unique ID assigned by PubMed to identify articles.

SupplMeshName	Supplementary conceptual term
SupplMeshName_Type	The type of supplementary conceptual term
SupplMeshName_UI	MeSH unique identifier for supplementary protocols and diseases

8. A08_ChemicalList

The chemical substances and registry number covered in each article. Registry Number refers to a code assigned by Chemical Abstracts Service to a specific chemical substance.

Table 8 Data Description for records of A08_ChemicalList

Column Name	Description
PMID	Unique ID assigned by PubMed to identify articles.
RegistryNumber	Unique Identifier to a specific chemical substance assigned by Chemical Abstracts Service
NameOfSubstance	The name of the specific chemical substance

9. A09_CommentsCorrectionsList

Reference information for each article, including the source, type, and PMID of the reference.

Table 9 Data Description for records of A09_CommentsCorrectionsList

Column Name	Description
PMID	Unique ID assigned by PubMed to identify articles.
RefSource	Reference sources
RefType	Reference type
RefPMID	The PMID of the reference
RefNote	Corrections to records with incorrect references

10. A10_DatabankList

The search number of the molecular sequence database that appears in the PubMed article. The search number can find the information of the corresponding chemical molecule from the established molecular sequence database, avoiding the use of lengthy molecular formulas and graphics in the article.

Table 10 Data Description for records of A10_DatabankList

Column Name	Description
PMID	Unique ID assigned by PubMed to identify articles.
DataBankName	Name of the molecular sequence database

11. A11_PersonalNameSubjectList

Individuals' names appear in <PersonalNameSubject> for citations that contain a biographical note or obituary, or are entirely about the life or work of an individual or individuals. Data is entered in the same format as A02_AuthorList.

12. A12_InvestigatorList

Each article corresponds to the NASA-funded principal investigator (PI) information, and they participated in the discussion and research of the article (but not necessarily the author).

Table 11 Data Description for records of A12_InvestigatorList

Column Name	Description
PMID	Unique ID assigned by PubMed to identify articles.
AffiliationInfo_Affiliation	The affiliation which the researchers belong to
AffiliationInfo_Identifier	Unique identifier of the affiliation

13. A13_AffiliationList

Extracted affiliation information.

Table 12 Data Description for records of A13_AffiliationList

Column Name	Description
PMID	Unique ID assigned by PubMed to identify articles.
Au_Order	Author order of the current authors in the author list of current articles.
Affiliation	Affiliation string.
Affiliation_order	For authors belonging to multiple affiliations, use numbers to identify different affiliations (in no particular order)

14. A14_ReferenceList

Reference information (we use the information in this table to generate the author's self-cited record during the disambiguation of strong features).

Table 13 Data Description for records of A14_ReferenceList

Column Name	Description
PMID	Unique ID assigned by PubMed to identify articles.
RefRank	Reference ranking
RefCitation	Reference sources (journals, year of publication, etc.)
RefArticleId	The PMID of the reference
RefIdType	Reference source database (e.g., PubMed)

15. B01_Descriptor

Specific information of Descriptor Name of each article (used for the classification of articles).

Table 14 Data Description for records of B01_Descriptor

Column Name	Description
DescriptorUI	The unique identifier of the descriptor
DescriptorName	The unique term of the descriptor
DescriptorClass	The corresponding classification of the descriptor (such as subject index, document type)
DateCreated	The date of the recording
DateEstablished	The date available for retrieval
Annotation	Text information corresponding to the descriptors
HistoryNote	History of text messages designed to help online searchers
NLMClassificationNumber	The NLM classification numbers assigned to the index terms
OnlineNote	Text messages designed to help online searchers
PublicMeSHNote	Change history of information
ConsiderAlso	Cross-reference, point to similar descriptors
PreviousIndex	Different usage in Descriptors and SCRs: In Descriptors: Free-text field referring to Descriptors or Descriptor/Qualifier combinations which were used to index the concept in the MEDLINE databases before the Descriptor was created.
TreeNumberList	One or more sets of elements in the descriptor or qualifier record

16. B02_Projects, B03_Link_Papers_Projects, B04_Link_ClinicalTrials_Projects, and B05_Link_Patents_Projects

Contains NIH funding information. Download link: <https://reporter.nih.gov/exporter>

17. C01_Papers

Simple table of A01_Articles.

Table 15 Data Description for records of C01_Papers

Column Name	Description
PMID	Unique ID assigned by PubMed to identify articles.
PubYear	Year the article was published.
AuthorNum	The number of authors of this article.
CitedCount	The number of times the article was cited (calculated from the C04_ReferenceList)

18. C02_Link_Papers_Authors

Simple table of A02_AuthorList.

Table 16 Data Description for records of C02_Link_Papers_Authors

Column Name	Description
PMID	Unique ID assigned by PubMed to identify articles.
PubYear	Year the article was published.
AuthorNum	The number of authors of this article.
Au_Order	Ranking of the current author's signature in the paper.
AID	Unique identifier for the author in our PKG.
BeginYear	The year this author published the first paper.

19. C03_Affiliations

Merging Vetle_Map and A13_AffiliationList information, including the parsed organization information, such as Zip code, Location, Country, etc.

Table 17 Data Description for records of C03_Affiliations

Column Name	Description
Id	Unique ID assigned to each record.
PMID	Unique ID assigned by PubMed to identify articles.
Au_Order	Author order of an article
AffiliationOrder	Affiliation order of one author, in case one author has multiple affiliations.
Affiliation	Affiliation text.
Department	Department parsed from Affiliation.
Institution	Institution parsed from Affiliation.
Email	Email parsed from Affiliation.
Zipcode	Zipcode parsed from Affiliation.
Location	Location parsed from Affiliation.
Country	Country parsed from Affiliation.
AID	Author disambiguation ID.
City	City parsed from Affiliation.
State	State parsed from Affiliation.
Vetle_Country	Country abbreviation.
Type	Affiliation type, such as COM, EDU, ORG, et.al.
Lat	Latitude of the Affiliation.
Lon	Longitude of the Affiliation.
Fips	Fips code.

20. C04_ReferenceList_Papers

The sources of data integration include PubMed's own citation data, NIH Open Citation Collection(run by iCite) and the OpenCitations Index of Crossref open DOI-to-DOI citations(run by OpenCitations).

Table 18 Data Description for records of C04_ReferenceList_Papers

Column Name	Description
Id	Unique ID assigned to each record.
PMID	Unique ID assigned by PubMed to identify articles.
RefPMID	PMIDs of references.

21. C05_PIs

Match the author AID and NIH Principal Investigator (PI) number PIID to generate correspondence table C05, including PIID, AID, CORE_PROJECT_NUMBER, PIID, Application_ID, etc.

Table 19 Data Description for records of C05_PIs

Column Name	Description
Id	Unique ID assigned to each record.
PI_ID	A unique identifier for each of the project Principal Investigators. Each PI in the RePORTER database has a unique identifier that is constant from project to project and year to year, but changes may be observed for investigators that have had multiple accounts in the past, particularly for those associated with contracts or sub-projects.
PMID	Unique ID assigned by PubMed to identify articles.
ProjectNumber	Project number.
SubProjectNumber	Sub project number.
AID	Author disambiguation ID.
Application_ID	A unique identifier of the project record in the ExPORTER database.

22. C06_Link_Papers_BioEntities

Entity information set extracted from document titles and abstracts using BERN2.

Table 20 Data Description for records of C06_Link_Papers_BioEntities

Column Name	Description
PMID	Unique ID assigned by PubMed to identify articles.
StartPosition	Start position of mention in an abstract.
EndPosition	End position of mention in an abstract.
Mention	Entity mentioned in an abstract.
Entityid(8 columns)	Normalized entity identifier, include mesh, mim, CL, cellosaurus, NCBITaxon, NCBIGene, CHEBI.
Type	Enumerated type of entity; values include species, disease, gene, drug, mutation, cell_line, cell_type, DNA, RNA.
is_neural_normalized	For diseases and chemicals, BERN2 use hybrid NEN models, which are a combination of both rule-based and neural network-based models. An entity that is not normalized by the rule-based model is then normalized by a neural network-based model.

23. C07_Authors

Disambiguated authors in PKG.

Table 21 Data Description for records of C07_Authors

Column Name	Description
AID	Unique author ID assigned by PKG.
BeginYear	The year this author published the first paper.
RecentYear	The year of the author's most recent publication.
PaperNum	The number of papers published by the author.
h_index	H-index of the author.

24. C10_Link_Papers_Journals

Journal details for the year the paper was published or the most recent year.

Table 24 Data Description for records of C10_Link_Papers_Journals

Column Name	Description
PMID	Unique ID assigned by PubMed to identify articles.
PubYear	Year the article was published.
Journal_ISSN	Unique identifier of the journal.
Journal_RecordYear	The year of the journal information for the row.
Journal_Title	Title of the journal.
Journal_SJR	The journal's SJR for the year.
Journal_SJR_Best_Quartile	The journal's best quartile according to its SJR score the year.
Journal_Hindex	The journal's H-index for the year.
Journal_Categories	The journal's category and corresponding quartile for the year.

25. C11_ClinicalTrials

Clinical trial studies from ClinicalTrials.gov. For more information, please refer to AACT.

26. C12_Link_Papers_Clinicaltrials

Citations between PubMed papers and clinical trial studies.

Table 25 Data Description for records of C12_Link_Papers_Clinicaltrials

Column Name	Description
PMID	Unique ID assigned by PubMed to identify articles.
nct_id	Unique ID assigned by ClinicalTrials.gov to identify clinical trial studies.
Type	Indicate the location of the reference or whether it is recorded in PubMed.

27. C13_Link_ClinicalTrials_BioEntities

Entity information set extracted from document titles and abstracts using BERN2.

Table 26 Data Description for records of C13_Link_ClinicalTrials_BioEntities

Column Name	Description
nct_id	Unique ID assigned by ClinicalTrials.gov to identify clinical trial studies.
StartPosition	Start position of mention in an abstract.
EndPosition	End position of mention in an abstract.
Mention	Entity mentioned in an abstract.
Entityid(7 columns)	Normalized entity identifier, include mesh, mim, CL, cellosaurus, NCBITaxon, NCBIGene, CHEBI.
Type	Enumerated type of entity; values include species, disease, gene, drug, mutation, cell_line, cell_type, DNA, RNA.
is_neural_normalized	For diseases and chemicals, BERN2 use hybrid NEN models, which are a combination of both rule-based and neural network-based models. An entity that is not normalized by the rule-based model is then normalized by a neural network-based model.

28. C14_Investigators

Investigators recorded in ClinicalTrials.gov, and their AID allocated in PKG.

Table 27 Data Description for records of C14_Investigators

Column Name	Description
nct_id	Unique ID assigned by ClinicalTrials.gov to identify clinical trial studies.
AID	Unique ID assigned by PKG to identify authors.
same_author_prob	The score of the linkage between the investigator and the AID.
affiliation_disambiguated	Institution id allocated by PKG base on the affiliation.

29. C15_Patents

Patents from USPTO. For more information, please refer to PatentsView.org.

30. C16_Link_Patents_Papers

Citations from patents to papers, it is the subset of pcs dataset run by Marx Matt. For more information, please refer to <https://zenodo.org/records/10215169>.

31. C17_Assignees

Disambiguated assignee data for granted patents from PatentsView, and organization disambiguation by OpenAlex's algorithm.

Table 28 Data Description for records of C17_Assignees

Column Name	Description
patent_id	Unique ID assigned by USPTO to identify patents.
assignee_sequence	Order in which assignee appears on the patent.
assignee_id	Unique PatentsView database assignee ID assigned by disambiguation algorithm.
disambig_assignee_organization	Organization name, if assignee is organization.
assignee_type	Classification of assignee: 2 - US Company or Corporation, 3 - Foreign Company or Corporation, 4 - US Individual, 5 - Foreign Individual, 6 - US Government, 7 - Foreign Government, 8 - Country Government, 9 - State Government (US). Note: A "1" appearing before any of these codes signifies part interest.
affiliation_disambiguated_OA	Unique ID and organization name assigned by OpenAlex's disambiguation algorithm.

32. C18_Link_Patents_BioEntities

Entity information set extracted from document titles and abstracts using BERN2.

Table 29 Data Description for records of C18_Link_Patents_BioEntities

Column Name	Description
PatentId	Unique ID assigned by USPTO to identify patents.
StartPosition	Start position of mention in an abstract.
EndPosition	End position of mention in an abstract.
Mention	Entity mentioned in an abstract.
Entityid(7 columns)	Normalized entity identifier, include mesh, mim, CL, cellosaurus, NCBITaxon, NCBIGene, CHEBI.
Type	Enumerated type of entity; values include species, disease, gene, drug, mutation, cell_line, cell_type, DNA, RNA.
is_neural_normalized	For diseases and chemicals, BERN2 use hybrid NEN models, which are a combination of both rule-based and neural network-based models. An entity that is not normalized by the rule-based model is then normalized by a neural network-based model.

33. C19_Inventors

Inventors recorded in USPTO and PatentsView. For more information, please refer to patentsview.org.

34. C21_Bioentity_Relationships

Biomedical entity relationships in papers.

Table 31 Data Description for records of C21_Bioentity_Relationships

Column Name	Description
PMID	Unique ID assigned by PubMed to identify articles.
Entity_id1/Entity_id2	The IDs of two biomedical entities.
relation_type	Type of the two biomedical entities.
relation_id	ID allocated by PKG to identify the relation.

35. C22_DatasetMethod

Using dictionary method, extract dataset and method entities from the article.

Table 32 Data Description for records of C22_DatasetMethod

Column Name	Description
PMID	Unique ID assigned by PubMed to identify articles.
Mention	Entity mentioned in an abstract.
Entity	Standard expression of entity.
Type	Identify whether this entity is a method or a dataset.

36. C23_BioEntities

Information of biomedical entities.

Table 33 Data Description for records of C23_BioEntities

Column Name	Description
EntityId	Normalized entity identifier.
Type	Enumerated type of entity; values include species, disease, gene, drug, mutation, cell_line, cell_type, DNA, RNA.
Mention	Text of the bioentity

37. C24_Link_Clinicaltrials_Patents

Citations between clinical trials and patents.

Table 34 Data Description for records of C24_Link_Clinicaltrials_Patents

Column Name	Description
nct_id	Unique ID assigned by ClinicalTrials.gov to identify clinical trial studies.
PatentId	Unique ID of patent in USPTO.
Type	Type of the citation; values include citing and cited.