

地名の情報学と データ駆動型研究の展開



北本 朝展（ROIS-DS人文学オープンデータ共同利用センター / 国立情報学研究所）

<http://codh.rois.ac.jp/>

歴史ビッグデータの統合解析

<http://codh.rois.ac.jp/historical-big-data/>



自然科学的
データ

人文
社会的
データ

天気

天候

地震

噴火

経済

人口

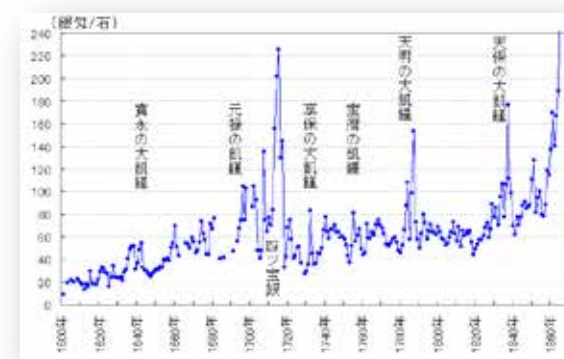
政治

文化

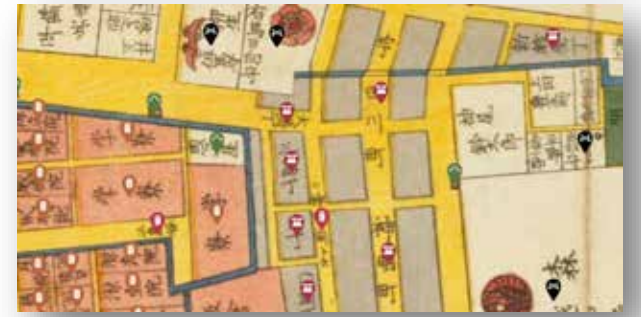
データ
構造化
ワーク
フロー

過去のビッグデータを統合解析するための基盤技術の研究

歴史ビッグ
データ基盤
(機械可読)



史料と地名のリンク



宇田川丁、三島丁・
神明丁、此分潰家多、
土蔵残所なし

固有表現認識

宇田川丁、三島丁・
神明丁、此分潰家多、
土蔵残所なし

御江戸大地震大破并出火類焼場等書上之写（みんなで翻刻）

1. 江戸マップに出現する地名を、**地名リソース（エンティティのデータベース）**として整備する。
2. 史料の文字列から**固有表現（地名）**を抽出する。
3. **特定のエンティティとリンク**することで、実世界と紐づける。

曖昧性解消

原資料表記	江戸マップID	江戸マップ表記
宇田川丁	4-358	宇田川町
三島丁	4-290	三島丁
神明丁	4-294	神明町

地名の情報学

空間情報としての
地名 = 地理情報
処理 (Geo)

+

テキストに出現
する地名 = 自然
言語処理 (NLP)

=

地名情報基盤
(Toponym
Information Platform)

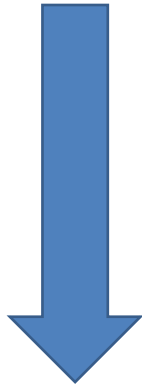
1. **GeoLOD** : 地名識別子の収集と共有を進める **ウェブサイト**
2. **Geoshape** : 地名識別子と関連付けた **地理形状データ** (境界など) の収集と共有を進める **ウェブサイト**
3. **GeoNLP** : テキストから地名を自動的に抽出しGeoLODの **地名識別子** と結合する **ソフトウェア**

地名識別子（地名ID）とは

1. 個々の地名をユニークな文字列（識別子 = ID）に対応させ、**すべての地名を識別子で区別可能**とする
2. **識別子のメタデータスキーマ**を定め、地名に関する各種の属性情報を管理する（緯度経度も含む）
3. 識別子は、**機械可読データを横断的に統合し、デジタル変革（DX）を推進する**要となる基盤データである
4. 識別子を維持・拡張していくには**安定した運用体制が必要**であり、**国家レベルで取り組む課題**である

地名識別子と疎結合性

私は**東京**に行く。



35.676666, 139.762222

地名を直接的に緯度経度に変換する方法は後から修正がしづらく、データ構造化と分析・可視化が一体化している（**密結合**）

私は**東京**に行く。



地名識別子：UoGwci



35.676666, 139.762222

地名を識別子にリンクし、識別子の属性の緯度経度を用いる方法は、分析・可視化の目的に応じて属性を変更できる（**疎結合**）

GeoLOD - 地名情報プラットフォーム

<https://geolod.ex.nii.ac.jp/>

The screenshot shows the GeoLOD website interface. At the top left is the GeoLOD logo and the text '地名情報を集約する地名情報処理システム'. A search bar contains the text '川崎' and a '検索' button. Below the search bar is a '結果一覧' section with a list of search results for '川崎' in various regions like '神奈川県川崎市' and '神奈川県横浜市'. A detailed information panel is open for the selected '川崎' entry, displaying the following data:

GeoLOD ID	MeYydo
地名	川崎
地名かな	
住所（現在）	神奈川県川崎市
緯度	35.530806
経度	139.703012
固有名称クラス	市区町村/政令指定都市
上位語	神奈川県/川崎市
説明	
異表記	
出典	1/川崎市役所/川崎市川崎区宮本町1/P34-14_14.xml
有効期限（始点）	1972-04-01
有効期限（終点）	
地名接頭辞	
地名接尾辞	市/

The background of the screenshot shows a map of Japan with numerous red location pins, and a specific pin for '川崎' is highlighted with a tooltip that says 'この付近を拡大'.

1. 登録した地名に**識別子（GeoLOD ID）**を付与し、アプリを越えて共有
2. 生成する地名語辞書は**GeoNLPの地名語辞書形式**とで活用できる形式
3. 有効期限など**歴史地名**にも対応

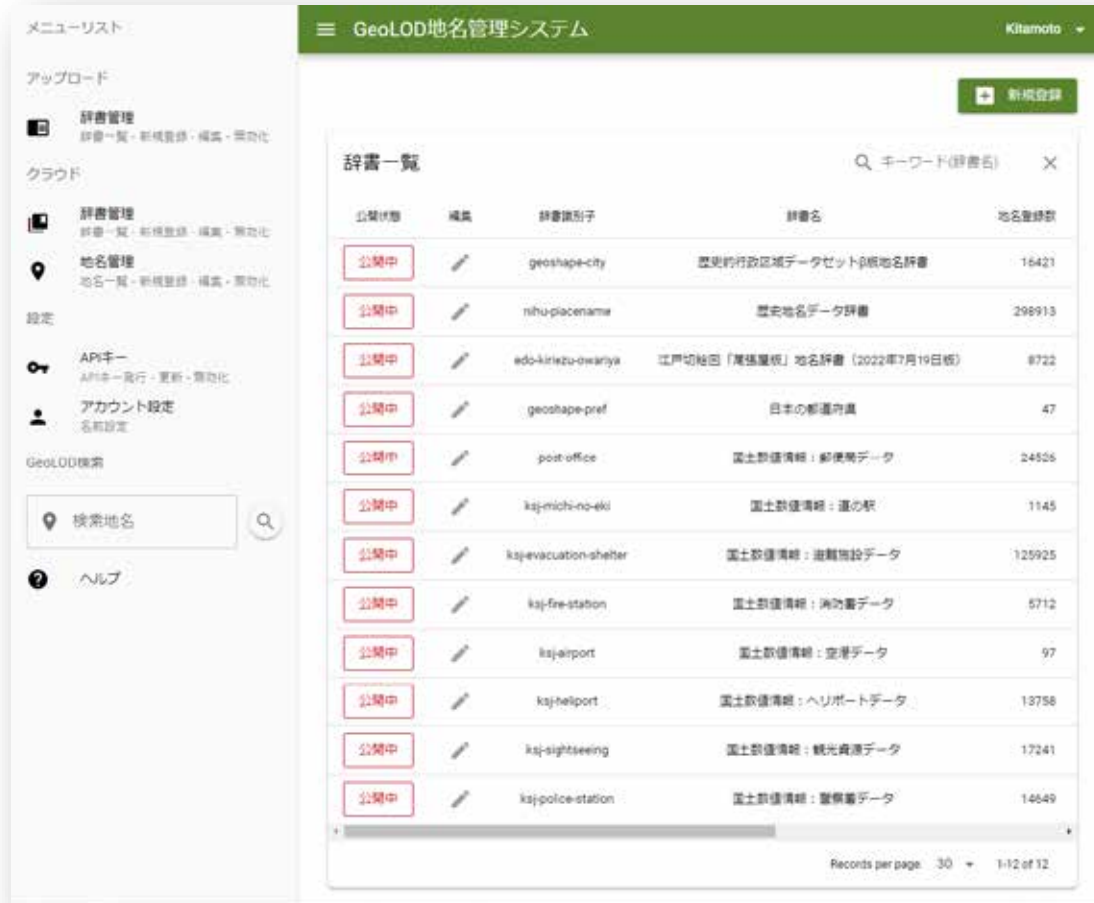
地名語辞書のスキーマ

<http://agora.ex.nii.ac.jp/GeoNLP/>

項目名	情報の種類	必須種別	説明
geolod_id	識別子	サーバ付与	GeoLOD内で一意のグローバル識別子
entry_id	識別子	必須	地名語辞書内で一意のローカル識別子
body	表記情報	必須	地名の原型
prefix	関係情報	推奨	接頭辞
suffix	関係情報	推奨	接尾辞
body_kana	表記情報	オプション	読み
ne_class	関係情報	必須	固有名クラス
hypernym	関係情報	推奨	上位語
latitude	属性情報	推奨	緯度（原則入力するが省略可）
longitude	属性情報	推奨	経度（原則入力するが省略可）
description	属性情報	オプション	説明
variant	属性情報	オプション	異表記
source	属性情報	オプション	出典（URL可）
valid_from	属性情報	オプション	有効期限（始点）
valid_to	属性情報	オプション	有効期限（終点）

GeoLOD地名管理システム

<https://geolod.ex.nii.ac.jp/admin/>



The screenshot shows the admin interface of the GeoLOD地名管理システム. The main content area displays a table titled '辞書一覧' (Dictionary List) with columns for '公開状態' (Publication Status), '編集' (Edit), '辞書識別子' (Dictionary ID), '辞書名' (Dictionary Name), and '地名登録数' (Number of Place Name Registrations). The table lists various dictionaries, all with a status of '公開中' (Public).

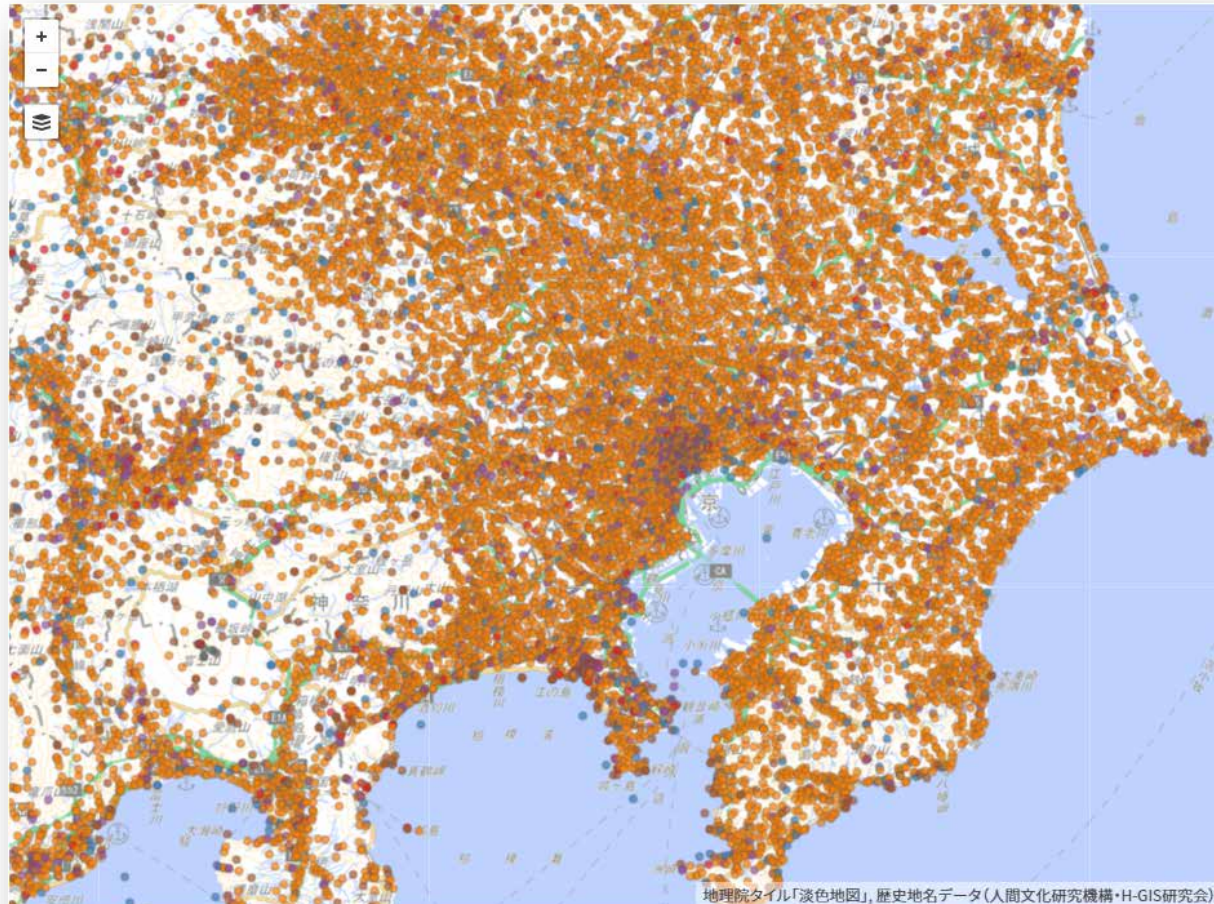
公開状態	編集	辞書識別子	辞書名	地名登録数
公開中	/	geoshape-city	歴史的行政区域データセット6版地名辞書	16421
公開中	/	nifu-placename	歴史地名データベース	298913
公開中	/	edo-kirizu-owariya	江戸切絵図「尾張屋敷」地名辞書 (2022年7月19日版)	8722
公開中	/	geoshape-pref	日本の都道府県	47
公開中	/	post-office	国土数値情報：郵便局データ	24526
公開中	/	kaj-michi-no-eki	国土数値情報：道の駅	1145
公開中	/	kaj-evacuation-shelter	国土数値情報：避難施設データ	125925
公開中	/	kaj-free-station	国土数値情報：海防署データ	5712
公開中	/	kaj-airport	国土数値情報：空港データ	97
公開中	/	kaj-heliport	国土数値情報：ヘリポートデータ	13758
公開中	/	kaj-sightseeing	国土数値情報：観光資源データ	17241
公開中	/	kaj-police-station	国土数値情報：警察署データ	14649

1. **アップロード辞書**：様々な地名集から収集した地名を、**オフライン**で編集し、**全部**を登録する

2. **クラウド辞書**：資料を読みながら足りない地名を、**オンライン**で編集・追加する

歴史地名マップ

<http://codh.rois.ac.jp/historical-gis/nihu-map/>



マーカーの色は歴史地名データの属性に対応し、**行政地名**、**建物**、**水部**、**地形**、**名所旧跡**、**その他** とします。また

1. 人間文化研究機構・H-GIS研究会が公開する「**歴史地名データ**」298,914件を活用
2. GeoLOD IDを新たに付与し、GeoNLP**地名語辞書形式**で公開
3. **バイナリベクトルタイル**の活用により、多数地点をウェブ地図に表示

江戸マップβ版

<http://codh.rois.ac.jp/edo-maps/>

番号	分類	現代語訳	翻刻	地図
2-001	施設	幸橋御門	幸橋御門	拡大図
2-002	施設	山下御門	山下御門	拡大図
2-003	施設	数寄屋橋御門	数寄屋橋御門	拡大図
2-004	施設	鍛冶橋御門	鍛冶橋御門	拡大図
2-005	施設	呉服橋御門	呉服橋御門	拡大図
2-006	地名	一石橋	一石橋	拡大図
2-007	地名	出橋	出橋	拡大図
2-008	町名	丸屋町	丸屋丁	拡大図

[2-296]
地名：磯辺大神宮（イソバ大神宮）
分類：寺社仏閣

29枚の江戸切絵
図から8722ヶ
所の地名を抽出
してデータベー
ス化

歴史的行政区域データセットβ版

<https://geoshape.ex.nii.ac.jp/city/>

行政区域は、多くの行政文書や統計データ、その他資料で頻出する、重要性が特に高い地名である



1. 国土数値情報の行政区域データセットを利用
2. 市区町村の連続性を判定し、独自の市区町村IDを付与
3. 市区町村の代表点などは公共施設データを参照して選定
4. 境界データと関連データの変遷を統合的に地図表示
5. データが連続していないため途中で抜けがあるのが問題

行政地名に関連するデータ

1. 行政区域名と境界の変遷に基づき、行政地名の連続性を判定し、識別子（市区町村ID）を網羅的に付与する
2. 全国地方公共団体コード：正式に定められた1968年以降しかない
3. 国土数値情報：1920年～2022年の行政地名と境界データを整備しているが、断続的で連続していない
4. 行政界変遷データベース（筑波大学 村山祐司研究室）：1889年～2006年の変化を連続的に追跡（エラーあり）
5. 『全訂 全国市町村名変遷総覧』（加除出版）など、出版社が著作権を有する資料も重要な存在

市区町村IDの付与



複数のデータセットを統合し、1920年以降の市区町村に網羅的なIDを付与

- 市区町村コード (XXXXXAYYYY) : 4161件
- 国土数値情報 (PPBQQQRRRR) : 12262件
- 筑波大データ (SSCTTTUUUU) : 358件
- 合計 : 16781件

- 1919年から1889年 (市制及町村制) への遡及
- 明治時代 (1888年以前) から江戸時代への遡及

出版社との協力



日本歴史地名大系
平凡社
巻冊数：50巻51冊
書籍版：1979-2004
現在、ジャパンナレッジで電子版を
公開中

1. 出版社が有する大規模・高品質データをどのように活用するか？
2. 日本歴史地名大系（平凡社）の「行政地名変遷表」を活用し、市区町村IDを江戸時代に接続できないか？
3. 一般社団法人「百科総合リサーチ・センター」が進める、行政地名変遷表のデジタルデータ公開に向けた作業に協力
4. 研究・教育に利用可能なデータを公開、ジャパンナレッジとも地名IDを共有予定

地名辞書の課題

1. 外国のカタカナ地名の辞書をどう作るか。表記揺れや言語の違いを扱えるか（例：キエフ vs. キーウ）
2. 日本国内のPOI（point of interest）の地名辞書をどう作るか。特にオープンデータの充実が課題
3. 自然地名、交通関係地名、公共施設名なども災害への利用では重要性が高い
4. 漢字地名の読みをどう付与するか。ニーズが大きい割にはデータが少ない
5. 地名の同一性や関係性（継承等）の収集、生活（SNS等）で使われる非公式の地名の収集など、課題は多数

GeoNLP – テキストジオタギング

<https://geonlp.ex.nii.ac.jp/pygeonlp/>

解析 クリア

特産品「明宝ハム」のメーカー、明宝特産物加工（郡上市明宝）の子会社、明宝マスターズが開設当初から運営に関わり、現在は指定管理者として駅全体の経営に関わる。駅長、マスターズ社長で明宝特産物加工の名畑和永専務（58）は「これだけ道の駅がある時代なので、他とは違うアイテム、役割を考え、差別化を図っていかねば」と生き残り策に苦心する。大手スーパーとは一線を画す、限定販売の高級明宝ハムなど各種商品、名物のケチャップをはじめ、昔ながらの漬物、プリンなど多彩な商品が並ぶ。

巨大な水車が目印の、恵那市山岡町の道の駅「おばあちゃん市・山岡」は来場者数が県の観光動態調査でも県内の観光施設で例年上位に入る。コロナ前の19年は58万9950人と高い水準だった。常に品切れさせない季節の野菜の直売が売りだ。先代の駅長がトヨタの「カンバン方式」さながらにジャストインタイムの仕入れ方法を発案。売り上げ規模は年2億円超と全国的には平均的だが、幹線道路から外れた立地ながらリピーターが絶えない。地域の活性化拠点という理想的な運営が光り、国内外からの視察が多い。

しかし、別の問題も浮かび上がっている。平成の大合併の名残で、一つの自治体に複数存在するケースがみられることだ。大垣市や多治見市など空白地もある一方で、郡上市は8カ所、高山市にも8カ所、下呂市には3カ所存在する。どこも似たつくりのため、同じ自治体の中で客を奪い合っている可能性もある。

1. Python版テキスト地名解析ツール
2. テキストから地名を抽出し、曖昧性を解消し、GeoLOD IDを自動付与し、GeoJSON形式で出力
3. Version 2.1では内部構造を変更し、ロジックを改良しやすくした
4. 今後は機械学習の導入を進める計画

Jageocoder - 日本の住所ジオコーダー

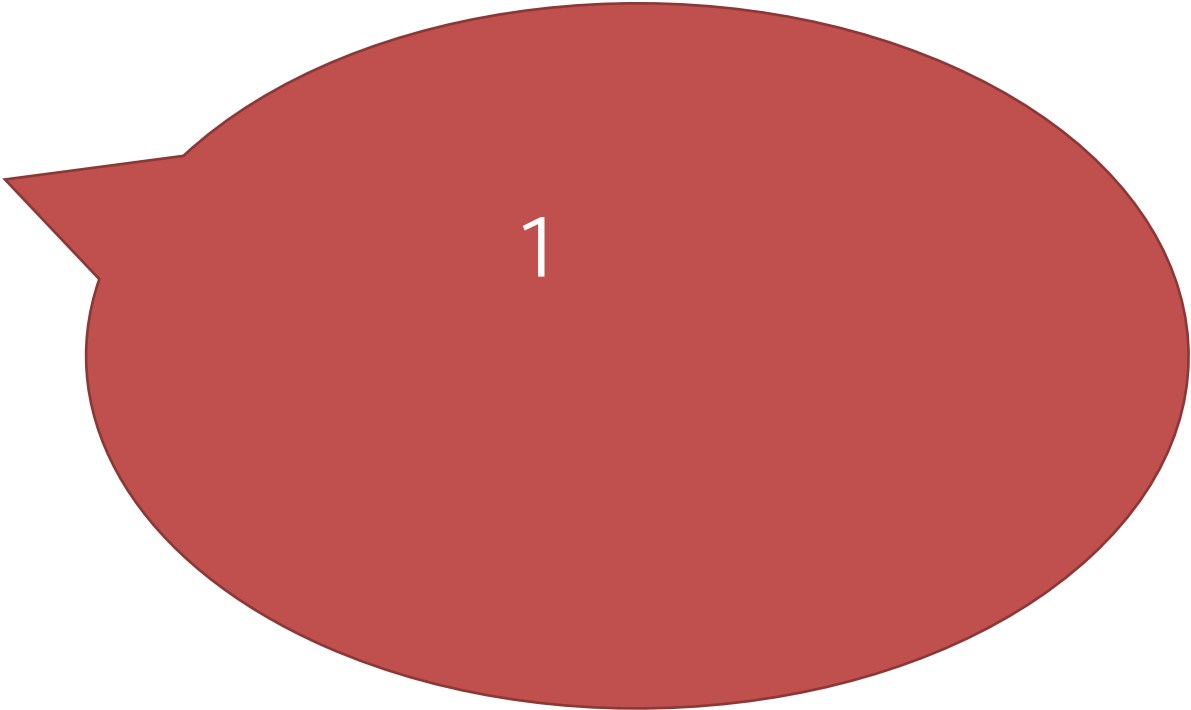
<https://geonlp.ex.nii.ac.jp/jageocoder/>



1. 情報試作室の相良毅氏が開発を進めるソフトウェア
2. Pythonで日本のアドレスを緯度経度に変換・正規化可能
3. 住居表示住所に加え地番住所の解析精度が大幅に向上
4. GeoNLPと連携させると、文章中の住所を自動的に抽出することも可能

歴史ビッグデータに必要な識別子

1. When → HuTime
2. Where → GeoLOD
3. Who → ?
4. What → ?



それぞれの課題に対して1つのサービスが生まれる。その全体をみんなで協力して運用する。

人名識別子や歴史イベント識別子など、他にも整備すべき識別子はたくさんある！

天文学における識別子の活用

<https://simbad.u-strasbg.fr/simbad/>

Query : M31

Basic data :
M 31 -- Galaxy
Other object types: LIN (), G (2006AJ,LEDA,...), * (AG,BD,...), GSO (2010A&A,[VV2006],...), AGN ([VV2000c],[VV2003c],...), *gam* (2FGL,3FGL,...), Rad (2C,DA,...), IR (IRAS,IRC,...), X (2MAXI,XSS), GiC (GIN), GiG (K79)

ICRS coord. (*ep=J2000*) : 00 42 44.330 +41 16 07.50 (Infrared) [] C 2006AJ....131.1163S
FK4 coord. (*ep=B1950 eq=1950*) : 00 40 00.095 +40 59 41.73 []
Gal coord. (*ep=J2000*) : 121.174329 -21.573309 []
Radial velocity / Redshift / cz : $V(\text{km/s}) -300.0 [4.0]$ / $z(\text{spectroscopic}) -0.001000 [0.000013]$ / $cz -299.85 [4.00]$
C 2012AJ....144....4M

Parallaxes (*mas*): 6.0 [14.1] E 1995GCTP..C.....0V
Morphological type: SA(s)b D 2013AJ....146...67B
Angular size (*arcmin*): 199.53 70.79 35 (Obj) D 2003A&A...412...45P
Fluxes (6) :
U 4.86 [0.03] D 2007A&JS..173..185G
B 4.36 [0.02] D 2007A&JS..173..185G
V 3.44 [0.03] D 2007A&JS..173..185G
J 2.094 [0.016] C 2006AJ....131.1163S
H 1.283 [0.017] C 2006AJ....131.1163S
K 0.984 [0.017] C 2006AJ....131.1163S

References (12047 between 1850 and 2023) (Total 12047)
Simbad bibliographic survey began in 1850 for stars (at least bright stars) and in 1983 for all other objects (outside the solar system).

Reference summaries :
from: 1850 to: \$currentYear
 or select by : (not exhaustive, [explanation here](#))

2022A&A...657A..150 [X ,2]
Astronomy and Astrophysics, volume 657A, 15-15 (2022/1-1)
A first estimate of the Milky Way dark matter halo spin.
OBREJA A., BUCK T. and MACCIO A.V.
Simbad objects: 6

2022A&A...657A..26M [X ,1]
Astronomy and Astrophysics, volume 657A, 26-26 (2022/1-1)
Discovery of four super-soft X-ray sources in XMM-Newton observations of the Large Magellanic Cloud.
MAITRA C. and HABERL F.
Simbad objects: 23

2022A&A...657A..41R [X ,1]
Astronomy and Astrophysics, volume 657A, 41-41 (2022/1-1)
The outermost stellar halo of NGC 5128 (Centaurus A): Radial structure.
REJKUBA M., HARRIS W.E., GREGGIO L., CRNOJEVIC D. and HARRIS G.L.H.
<Available at CDS ([J/A+A/657/A41](#)): fields.dat catalog.dat>
Simbad objects: 7
Status at CDS: *Tables of objects will be appraised for possible ingestion in SIMBAD.*

すべての天体に識別子が付与され、
各種の情報が統合されているため、
例えば「ある天体に関する論文」な
どが簡単に検索できる

データ駆動型研究の未来

1. 多くの研究者が標準化された地名識別子を活用すれば、「ある地名に関する研究成果」なども一覧できる
2. 多くのグループが共通の地名識別子を用いてデータを作成すれば、データの再利用性が画期的に向上する
3. 分野ごとに地名辞書を増強していけば、他の分野もその成果の恩恵を受ける
4. 過去から現在、未来にわたる時間的な変遷を追跡できれば、未来に役立つ知識が得られる

もっと詳しく



← 本日の発表資料

1. 北本 朝展, 村田 健史, "歴史的行政区域データセットβ版をはじめとする地名情報基盤の構築と歴史ビッグデータへの活用", 情報処理学会技術報告, Vol. 2020-CH-124, No. 1, pp. 1-8, 2020
2. 北本 朝展, 鈴木 親彦, 寺尾 承子, 堀井 美里, 堀井 洋, "地理的史料を対象とした歴史地名の構造化と統合に基づく江戸ビッグデータの構築", 人文科学とコンピュータシンポジウム じんもんこん2020論文集, pp. 171-178, 2020
3. 北本 朝展, "地名情報基盤GeoLODによる地名識別子の収集・共有・活用と歴史ビッグデータ研究", 人文科学とコンピュータシンポジウム じんもんこん2022論文集, pp. 7-14, 2022



ROIS-DS人文学オープンデータ共同
利用センター

<http://codh.rois.ac.jp/>