

国立国語研究所の言語資源と オープンデータ・オープンサイエンス

小木曽智信 2019/10/24 NII



自己紹介

- 小木曽 智信(おぎそとしのぶ) togiso@ninjal.ac.jp
 - 人間文化研究機構 国立国語研究所 言語変化研究 領域代表 (コーパス開発センター兼任)
 - 共同研究プロジェクト「<u>通時コーパスの構築と日本語</u> 史研究の新展開」プロジェクトリーダー(2016年~)
 - 専門は日本語学、自然言語処理
 - 『太陽コーパス』『現代日本語書き言葉均衡コーパス』『日本語歴史コーパス』など書き言葉のコーパス 構築に携わる



国語研と オープンサイエンス・オープンデータ



オープンサイエンス・オープンデータを 基盤に

- ・オープンなデータ
 - 調査収集したデータは公開を原則とする
 - − 所定の手続きで誰でもアクセス可能に※必ずしも無償・無制限であることまでは意味しない
- 研究方法もオープンに、検証可能に
 - 実験など主観を排した方法で
 - 研究に用いた中間データやスクリプト類もできるだけ 公開

第4期の国語研の基本方針(案)より

コーパスとアーカイブを核に

- これまでに高い評価を受け、今期の「多様な言語資源に基づく総合的日本語研究の開拓」でも多数構築中のコーパスを中心とする言語資源を更新・拡張
- 危機言語データ等のアーカイブに取り組み、 オープンサイエンスを支える

第4期の国語研の基本方針(案)より

次期、オープンサイエンスを謳う理由

現在はまだできていないから

- データの囲い込み(個人作成・収録データ)
- エビデンスが見えない文法性判断に基づく議論
- 論文の元となった(中間)データの非公開・・・

それでも、それに向けた取り組みは始めている



国語研のコーパスとオープンデータ



国立国語研究所のコーパス

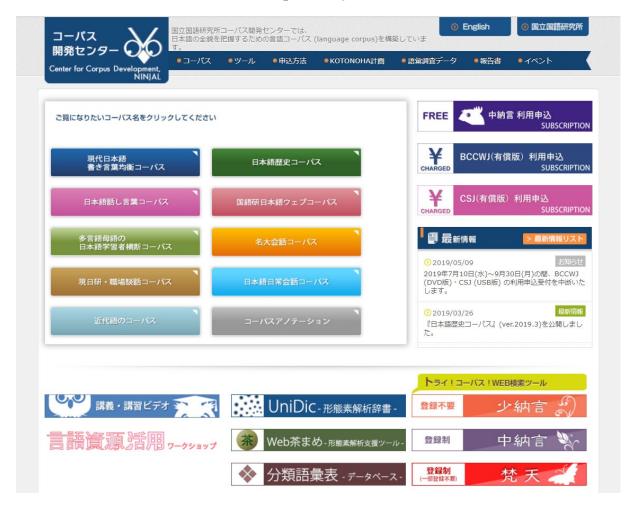
コーパスとは:言語を分析するための基礎資料として,書き言葉や話し言葉の資料を体系的に収集し,研究用の情報を付与した大規模なデータベース

- 『日本語話し言葉コーパス』(CSJ) 2004年
- ・『現代日本語書き言葉均衡コーパス』(BCCWJ) 2011年
- 『日本語歴史コーパス』(CHJ) 2013年~
- ・『国語研日本語ウェブコーパス』 2016年
- ・『多言語母語の日本語学習者横断コーパス』2016年~ ※現在は日常会話、方言などのコーパスも構築中





コーパス開発センター



研究インフラとしてのコーパス

- •「現代日本語書き言葉均衡コーパス」
 - 登録ユーザ数: 約20,000人
 - 年間クエリ数:約50万件/年
 - 利用した論文数:約70本/年
- •「日本語歴史コーパス」
 - 登録ユーザ数:約10,000人
 - 年間クエリ数:約26万件/年
 - 利用した論文数:約50本/年(※予稿集含む)



コーパス検索アプリケーション「中納言」

- 人文系の研究者に利用しやすい形で提供
 - オンライン
 - -無料(要登録)





コーパス構築のコスト

「現代日本語書き言葉均衡コーパス」(BCCWJ) 2006~2010年(研究代表者:前川喜久雄)

- ・ 約1億語の現代語書き言葉のコーパス
 - 書籍等からのバランスをとったサンプリング
 - 紙の文献の電子化
 - 単語情報の付与
- 予算:
 - 科研費 特定領域研究(総額約8億円)
 - 十国語研運営費交付金



コーパス構築のコスト2

- 「日本語歴史コーパス」(CHJ) 2013年~現在(研究代表者:小木曽)
- 奈良時代から明治・大正時代までの日本語
 - 全本文に単語情報の付与
 - 外部の画像等にリンク
 - ※2013年以前に構築されたデータを継承して含む
- 予算:
 - 国語研運営費交付金より年間3000万弱 +科研費 基盤(A) 年間1000万弱



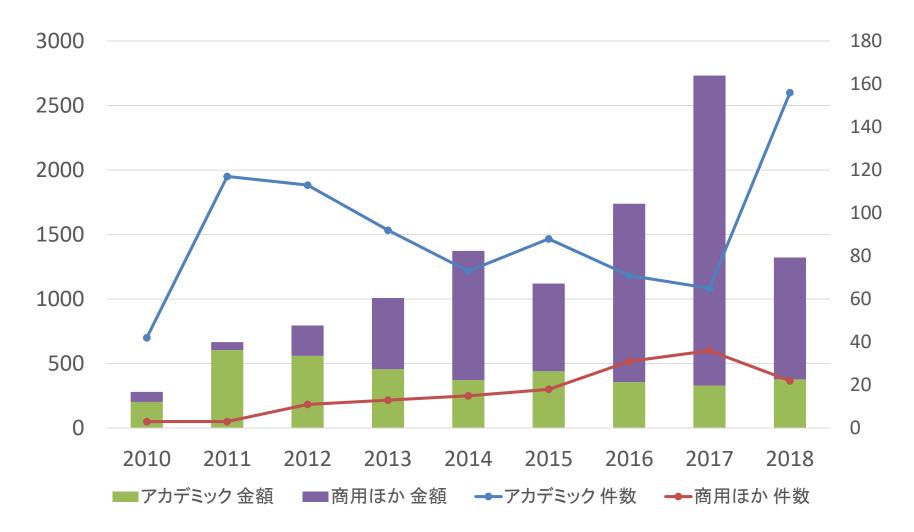
コーパス公開のコスト

サーバー代等: 概ね年間1000万円程度 (十人件費、電気代 etc.)

新しいコーパスを構築する際は、その新規性をもって予算獲得ができるが、インフラとして定着したコーパスの公開・維持費は外部資金では困難



コーパスからの収益





コーパスとオープンデータ

- コーパスは研究の副産物ではなく、その構築自体が 目的でもある
- コーパスは他者の著作物ではなく、構築者自らの(編集)著作物である
- コーパスは自己収入を生み出す経済的価値を持つ
- > 完全なオープン化は困難
- ▶ 所定の手続きで誰でもアクセス可能にするが、※必ずしも無償・無制限であることまでは意味しない



コーパスとオープンデータ

コーパス本体はオープンデータ化は難しいが、コーパスに対するアノテーションは公開できる

- アノテーション: 追加情報の付与
 - ・コーパス中の単語や文、発話・・・に
 - 意味情報、統語情報、メタ情報・・・を付与
- コーパスに依拠したオープンサイエンスの可



『日本語歴史コーパス』を例に

コーパスとオープンサイエンス



これまでの日本語研究では、コーパスを利用した場合でも、コーパス中の用例を分類したり分析をおこなったりした資料は個人のもとで管理され公開されることはなかった。研究が論文にまとめられると、その基礎となったデータは再利用されず死蔵され、場合によっては破棄される場合も少なくない。



しかし、コーパスに依拠したデジタルデータは、 インターネット上で共有することは比較的容 易である。基礎データの共有が可能になれば、 当該の研究が検証可能になるだけでなく、そ の成果を再利用して他の研究者が新たな研 究を行うことも可能になる。コーパスを対象と した研究データの多くは、コーパスに対するア ノテーションとしてまとめられる。



オープンサイエンスに向けたコーパス 検索ツール側の仕組み

- 1. 検索条件式とその共有
- 2. 用例へのパーマリンク
- 3. 用例のユニークIDとアノテーション
- 4. アノテーションデータの共有
- 5. アノテーション共有環境の構築



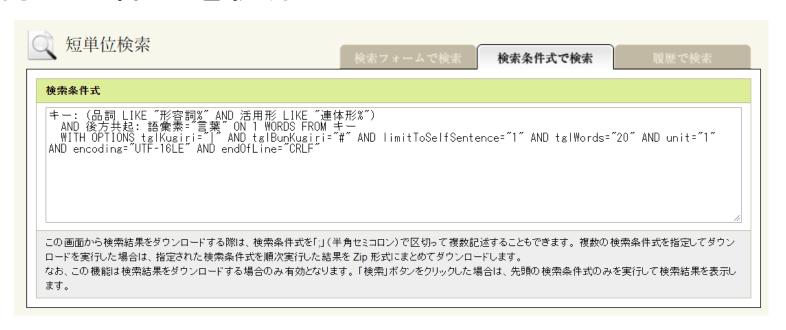
1. 検索条件式

・「中納言」で"「言葉」という語の直前に来る形容詞の連体形を検索"



1. 検索条件式

・「中納言」で"「言葉」という語の直前に来る形容 詞の連体形を検索"



・ 研究の再現性、用例の共有



1. 検索条件式

• 「検索条件式」を使うことで、中納言ユーザーなら誰でも、同じ検索を行うことができる

キー: (品詞 LIKE "形容詞%" AND 活用形 LIKE "連体形%") AND 後方共起: 語彙素 = "言葉" ON 1 WORDS FROM キーWITH OPTIONS unit="1"

- 研究の再現性のために論文などで使用した 検索条件式を明記するとよい
- →研究データの共有、検証可能性



2. 用例へのパーマリンク

• 位置情報をもとにした中納言上の用例への リンク

https://chunagon.ninjal.ac.jp/chj/permalink?un
it=short&position=20-源氏1010_00012,134430

- 論文や資料のPDFに用例へのリンクを埋め込み!
- QRコードにも! SNSでも!



2. 用例へのパーマリンク





3. 用例のユニークID

『日本語歴史コーパス』の位置情報

- サンプルIDと開始位置 (=先頭からの何文字目か)
 - 20-源氏1010_00012, 134430
- 位置情報で用例が一意に定まる=<u>用例の</u> マイナンバー

サンブル ロ \$	開 始 位 置	連番 ◆	コ ア •	前文脈◆	‡ - •	後文脈 ◆	語 形 ◆	品 詞 ◆
20-源氏 1010_00012	134430	77920	1	ぞ つと さぶらひ ける 。# 前栽 の 花 いろ いろ 咲き 乱れ 、 おもしろき 夕暮 に、 海 見 やら	ත ත	廊 に 出でたまひ で 、 たたずみ たまふ 御 さま の ゆゆしう きょら な る こと 、 所がら は まして こ	ル	助動詞
				うまつるをうれしきことして、四五人ばかりぞつ とさぶらひける。# 前栽の花いろいろ咲き乱れ、おもしろき夕暮に、海見やら	5 5	廊に出でたまひて、たたずみたまふ 御さまのゆゆしうきょらなるこ		



3.用例のユニークID

用例集としての位置情報

- 用例の位置情報を集めて並べればそれだけで 貴重な用例集になるかも!
 - あとで再利用できる、検証できる
 - 「中納言」経由で現代語訳や原文画像が参照できる

位置情報を用いたアノテーション

- 位置情報に、用例の情報を追加すれば立派なアノテーションデータ!
 - 20-源氏1010_00012, 134430, 自発



4. アノテーションデータの共有

小木曽(2019)<u>『日本語歴史コーパス』への追加情報の付与と共有―中古和文の「る」「らる」を例に―</u>日本語学会2019年度春季大会予稿集

サンプル ID	開始位置	連番	本文種別	語形	用法
20-伊勢0920_00001	1730	1070	歌	ル	可能
20-伊勢0920_00001	29880	19170		ラル	<u>受身</u>
20-伊勢0920_00001	68940	44330		ラル	<u>受身</u>
20-伊勢0920_00001	73050	47020	歌	ル	<u>自発</u>



4.アノテーションデータの共有

https://researchmap.jp/mu1dor8so-12361/

🍑 資料公開

資料公開 >> コンテンツ詳細

 タイトル
 CHJ中古「る」「らる」用法分類アノテーションデータ

 カテゴリ
 研究データ

 概要
 『日本語歴史コーパス』中古「る」「らる」用法分類アノテーションデータ ver.0.5

 ダウンロード
 CHJ中古れる・られるv05.csv (25)

 利用条件
 CC-BY

 記入者:togiso

-覧へ戻る



4.アノテーションデータの共有

- データの配布と引用
 - 小木曽(2019)「CHJ 中古「る」「らる」用法分類アノ テーションデータ」

https://researchmap.jp/mu1dor8so-12361/

- オープンデータとして「公開」しよう!
- アノテーションを「引用」しよう!!
 - → コーパスをみんなで「育て」よう!!!



4.アノテーションデータの共有

「……ともに「コーパスを育てていく」作業に参加 してほしい。こうして作られ、公開されたアノ テーションデータは他の研究者に役立つ共有 の財産となるものであるから、このようなデータ の作成・公開が、研究の業績として正当な評価 を得られるようになることも願っている。」 小木曽(2019)



5.アノテーション共有環境の構築

- ユーザーがもっと容易にアノテーションを 付与し、共有し、引用し、評価しできる環境を アプリケーション上で実現したい・・・
- ユーザーが多い「中納言」の機能拡張で・・・



科研·挑戦的研究(開拓) 19H05477

「日本語コーパスに対する情報付与を核とした オープンサイエンス推進環境の構築」

- 研究代表者: 小木曽智信
- 2019~2021年度
- 25,740千円





科研·挑戦的研究(開拓) 19H05477

- •「中納言」に、利用者が新たな情報を任意の場所に付加するアノテーション機能を追加
- 付加した情報を他のユーザーと共有できるシステムを構築
- 公開されたアノテーションデータを引用・評価できる仕組みを提供
- ▶この環境の構築とアノテーションの実践により、研究データの公開と共有を促し、書き言葉テキストを中心としたオープンサイエンスの基盤とすることを目指す。



まとめ

- 国語研のコーパスとその検索環境はすでに 日本語研究のインフラとして機能している。
- コーパスの完全なオープンデータ化は困難だが、コーパスとアノテーションの共有を通してオープンサイエンスを推進する基盤となしうる。
- 検索環境を拡張し、言語研究のオープンサイ エンス推進環境の整備をすすめ、今後につな げたい。

