

# J-UniMorph: Japanese Morphological Annotation through the Universal Feature Schema

Kosuke Matsuzaki<sup>1</sup>, Masaya Taniguchi<sup>2</sup>, Kentaro Inui<sup>3,1,2</sup>, Keisuke Sakaguchi<sup>1,2</sup>  
(<sup>1</sup> Tohoku University, <sup>2</sup> RIKEN, <sup>3</sup> MBZUAI)

- ✉ matsuzaki.kosuke.r7@dc.tohoku.ac.jp
- 𝕏 @Matsuzaki\_NLP
- 🌐 <https://matsukosuke.github.io/>

SIGMORPHON2024@Mexico City, June 20th



# Overview

- **Introduction**

- Prior Japanese dataset in UniMorph

lemma	inflected form	label
eat	ate	V;PST

Illustration of UniMorph

- **J-UniMorph: Japanese Adaptation**

- Creation process of J-UniMorph
- Comparison with the prior dataset

- **Application: J-UniMorph Visualizer**

- Aiming to support Japanese learners  
in semantic understanding and  
morphological usage

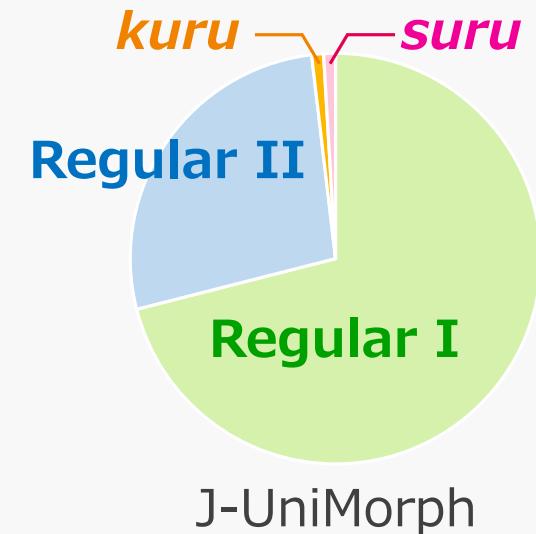


 unimorph/jpn

# Contributions to UniMorph [Sylak-Glassman'16]

- Increased Japanese Inflection Patterns and Combinations

- The number of inflection patterns becomes about 10x larger than the prior dataset

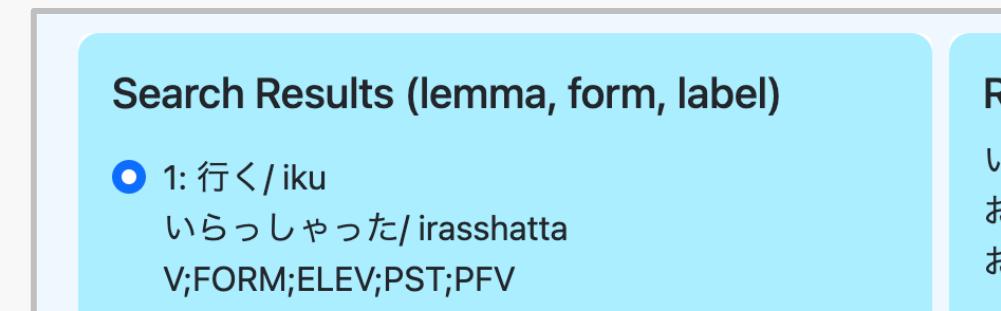


- Included High-Frequency Forms

- Average Google Search Hits became higher

- Improved Data Accessibility

- J-UniMorph Visualizer made it easier to explore and analyze



# Prior Japanese Dataset in UniMorph

- Released for the SIGMORPHON–UniMorph 2023 Shared Task 0
  - Automatically extracted from Wiktionary [Goldman+’23]
  - We call this **Wiktionary Edition**
- Basic statistics of the Wiktionary Edition:
  - 1. About 70% of verbs have the **same inflection pattern** †
    - These inflections can be replaced by one seed verb (lemma)
  - 2. There are an **average of 12 entries** for each verb
    - Japanese verbs have a more diverse range of inflections

---

†Details on the next slide

# Derived Verbs vs. Common Verbs

- Regular I verbs

- e.g., *kak-u* (write-V), *hashir-u* (run-V)

- Regular II verbs

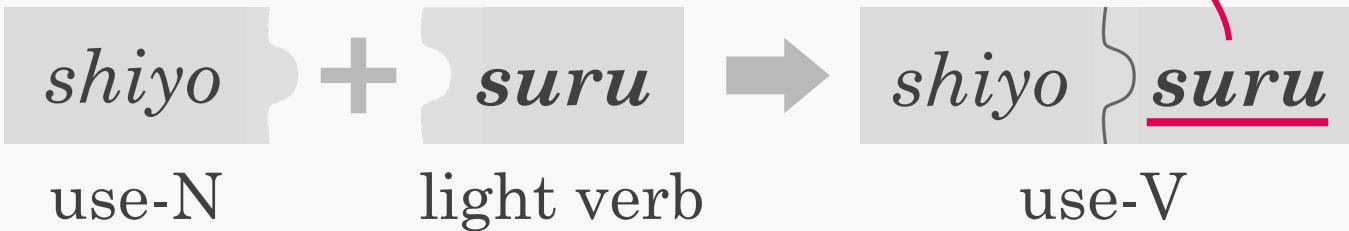
- e.g., *tabe-ru* (eat-V), *mi-ru* (see-V)

- Irregular verbs

- Only *kuru* (come-V), *suru* (do-V)

Derived verbs

- e.g.,



# Creation Process of J-UniMorph

## 1. Inflection Generation

- Generating the forms with verb inflection tool (semi-automatically)



## 2. Label Annotation Mapping

- Annotating UniMorph labels to each inflection pattern

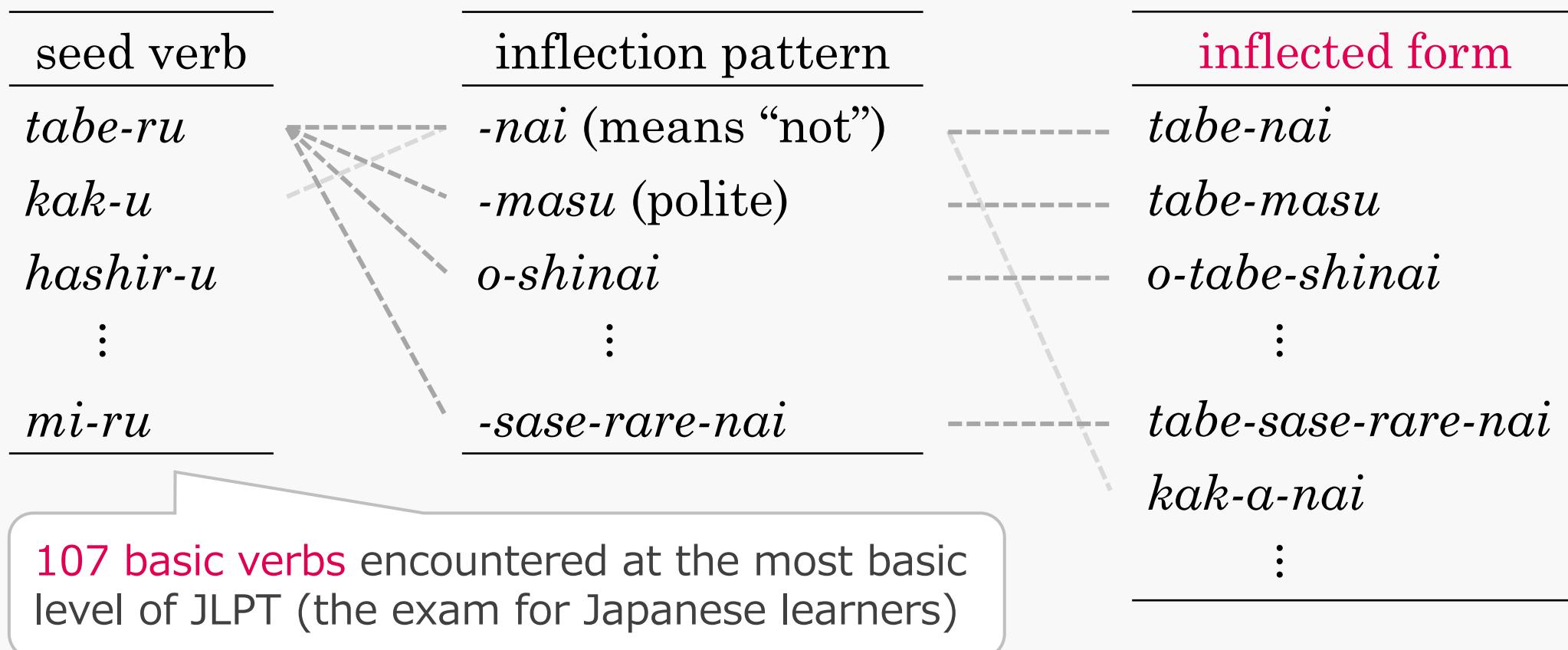


## 3. Filtering

- Remove incorrect or infrequent forms based on frequency

# Dataset creation Step 1: Inflection Generation

- Generate inflected forms with the rule-based verb inflection tool



# Dataset creation Step 2: Label Mapping

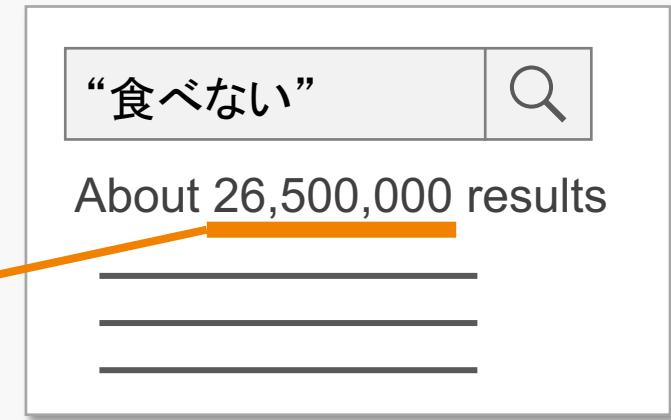
- Annotate the UniMorph labels to each inflection pattern (affixes)
  - References:
    - Definitions of UniMorph feature schema [Sylak-Glassman'16]
    - Books on Japanese grammar [Nitta'07, '09a, '09b; Hirabayashi+'88; Kato+'89; Kamiya'01; Takami'11]

inflection pattern	UniMorph feature label
<i>-nai</i>	V;PRS;IPFV;NEG
<i>-masu</i>	V;PRS;IPFV;POL;FOREG
<i>o-shinai</i>	V;PRS;IPFV;FORM;HUMB;NEG
:	:
<i>-sase-rare-nai</i>	V;PRS;IPFV;CAUS;PASS;NEG

# Dataset creation Step 3: Filtering

- Remove unnatural inflected forms
  - Using the number of hits from Google's exact match search
  - Remove forms with fewer than 10 hits

Inflected form	Romanization	Hits
食べない	<i>tabe-nai</i>	26,500,000
食べます	<i>tabe-masu</i>	22,600,000
お食べしない	<i>o-tabeshinai</i>	6
:	:	:
見させられない	<i>mi-sase-rare-nai</i>	137,000



← remove ( $\leq 10$ )

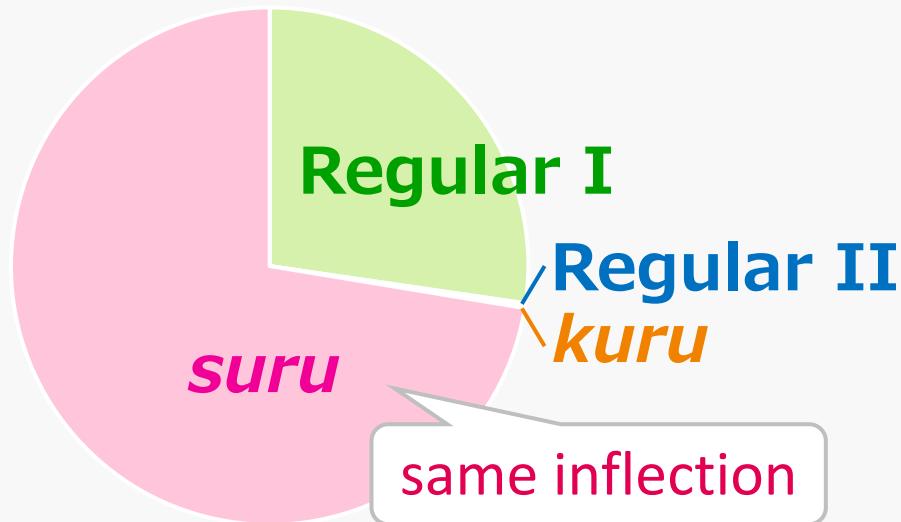
# Comparison (1/3): Average Entries per Word

- J-UniMorph has various inflection patterns and combinations
  - About 10x larger than the Wiktionary Edition despite similar entries

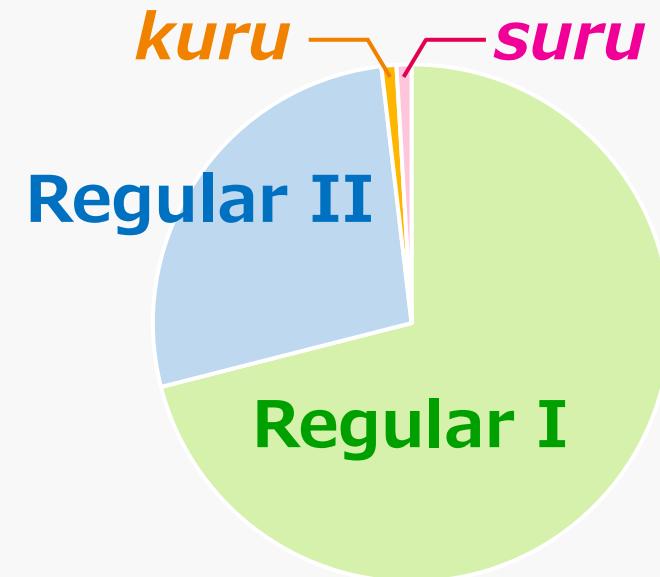
	Wiktionary Edition	J-UniMorph	Ours
Total number of entries	12,000	12,687	
Number of entries <b>per word</b>	12.0	118.6	

# Comparison (2/3): Verb Statistics

- Wiktionary Edition:
  - About 70% are the verbs which end with the light verb “*suru*”
  - These inflections can be replaced by the single verb “*suru*”



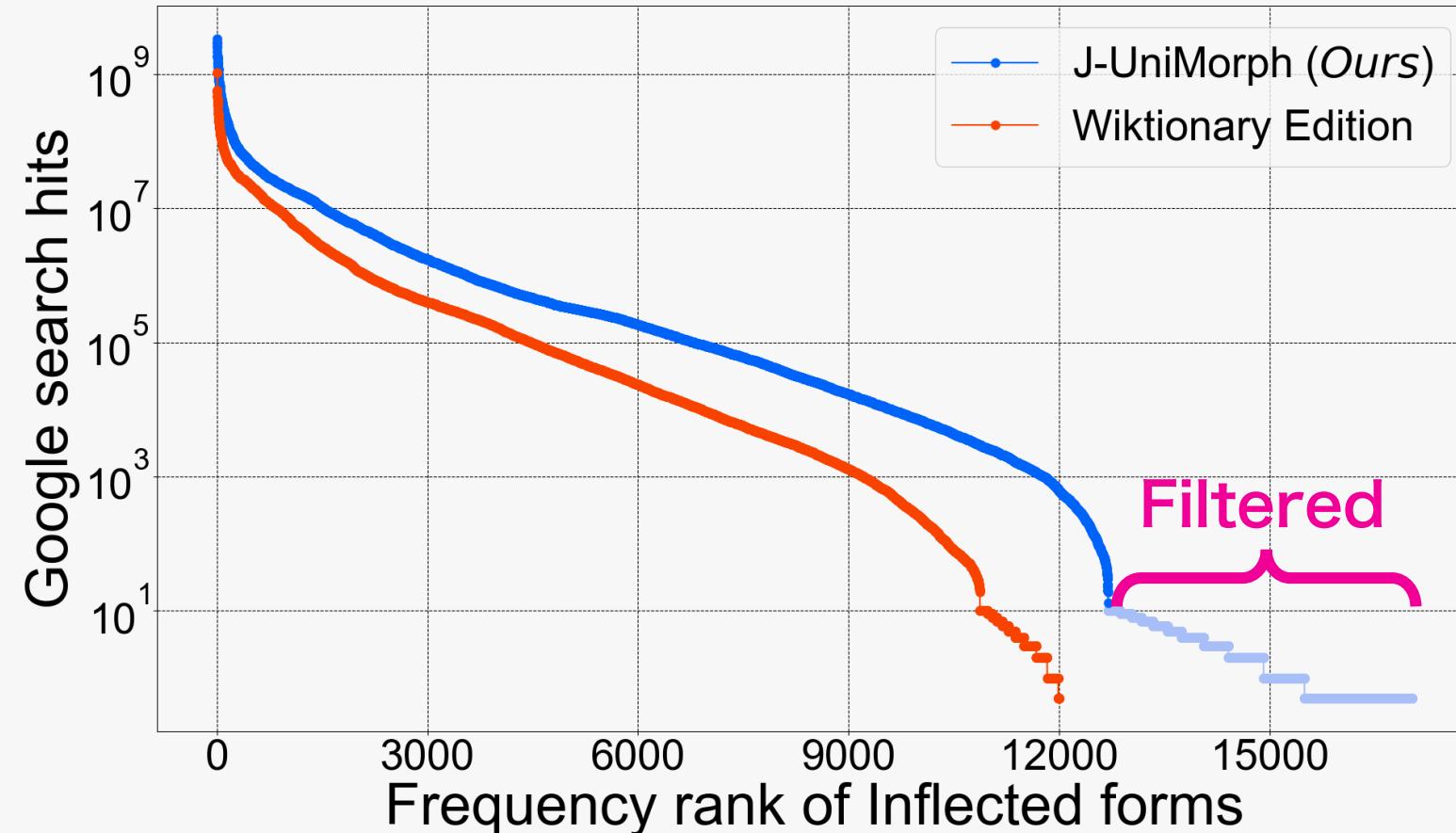
Wiktionary Edition  
(1,000 verbs / 12,000 entries)



J-UniMorph  
(107 verbs / 12,687 entries)

# Comparison (3/3): Hits of each Inflected Form

- The entries in J-UniMorph have higher frequency forms



# Application: J-UniMorph Visualizer

- Creating J-UniMorph
  - Dataset creation process
  - Comparison with the Wiktionary Edition



- **J-UniMorph Visualizer**
  - Aims to support Japanese learners in semantic understanding and morphological usage

↑ Visualizer

# Visualizer (1/3): Semantic Understanding



You can analyze the morphological meaning of inflected forms

Input 1: an inflected form

Output 1: lemma, morphological meaning, and related words

lemma: taberu  
meaning: V;PST;PFV;CAUS  
related: tabesaseta (67.11)  
tabesashita (38.75)

$$\text{FrequencyIndex} = 100 \times \frac{\log_{10}(\text{hits} + 1)}{\log_{10}(\text{max\_hits} + 1)}$$

“lemma” and “meaning” can be displayed in multiple instances

# Visualizer (2/3): Morphological Usage

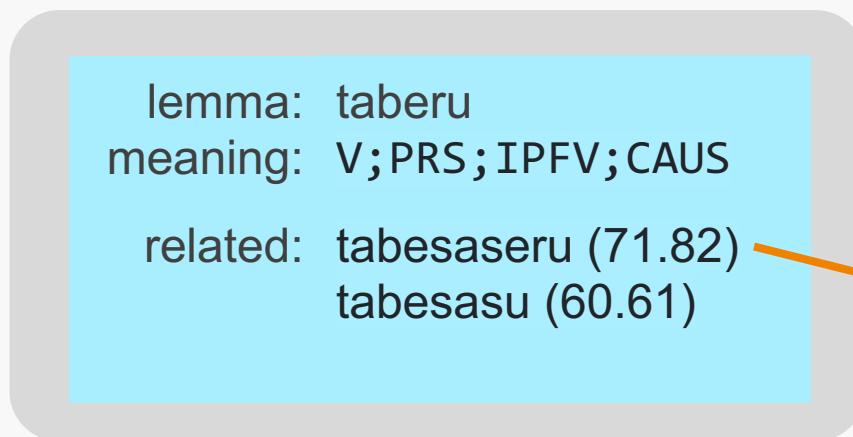


You can display other inflected forms with different labels

Input 2: switch the checkboxes to display other forms

tabesaseta	<input checked="" type="checkbox"/> V (Verb)	<input checked="" type="checkbox"/> PRS (Present)	<input checked="" type="checkbox"/> IPFV (Imperfective)
	<input checked="" type="checkbox"/> CAUS (Causative)	<input type="checkbox"/> PST (Past)	<input type="checkbox"/> PFV (Perfective)

Output 2: words which have new meanings you clicked



$$\text{FrequencyIndex} = 100 \times \frac{\log_{10}(\text{hits} + 1)}{\log_{10}(\text{max\_hits} + 1)}$$

# Visualizer (3/3): Screenshots



## 1: See lemma and meaning

irasshatta

Search Examples: 走りません, いらっしゃった, 食べられる, hashirimasan, irasshatta, taberareru

Input 1

Search Results (lemma, form, label)

- 1: 行く / iku  
いらっしゃった / irasshatta  
V;FORM;ELEV;PST;PFV
- 2: 居る / iru  
いらっしゃった / irasshatta  
V;FORM;ELEV;PST;PFV
- 3: 来る / kuru  
いらっしゃった / irasshatta  
V;FORM;ELEV;PST;PFV

Related Words

いらっしゃった / irasshatta (71.26)  
おいでになった / oideninatta (58.06)  
お行きになった / oikininatta (40.83)

Output 1

Part of Speech: V (Verb)

Tense: PRS+IPFV (Present, Imperfective) (unchecked), PST+PFV (Past, Perfective) (checked)

Honorifics: POL (Polite) (unchecked), FORM (Formal) (checked), HUMB (Humbling) (unchecked)

Mood: OPT+1 (Optative, Subjective) (unchecked), IMP (Imperative) (unchecked), INTEN (Intensive) (unchecked), PERM (Permissive) (unchecked)

Others: NEG (Negative) (checked), CAUS (Causative) (unchecked)

## 2: See other inflected forms

Part of Speech: V (Verb)

Tense: PRS+IPFV (Present, Imperfective) (checked), PST+PFV (Past, Perfective) (unchecked)

Honorifics: POL (Polite) (checked), FORM (Formal) (checked), HUMB (Humbling) (unchecked)

Mood: OPT+1 (Optative, Subjective) (unchecked), IMP (Imperative) (unchecked), INTEN (Intensive) (unchecked), PERM (Permissive) (unchecked)

Others: NEG (Negative) (checked), CAUS (Causative) (unchecked)

Uncheck

Check

Input 2

<Words with "行く V;PRS;POL;FOREG;NEG;FORM;ELEV;IPFV">

いらっしゃいません / irasshaimasen (76.63)  
おいでになりません / oideninarimasen (52.17)  
お行きになりません / oikininarimasen (32.76)

Output 2

# Extended dataset over UniMorph Schema



- We created the additional dataset for J-UniMorph Visualizer
  - <https://github.com/cl-tohoku/J-UniMorph>
- It includes (without Filtering):
  - Lemma (Kanji & Hiragana & Romanized form)
  - Inflected form (Kanji & Hiragana & Romanized form)
  - UniMorph Feature Labels
  - Google Search Hits

Lemma	Inflected form	Feature label	Hits
書く / かく / kaku	書かない / かかない / kakanai	V;PRS;IPFV;NEG	16000000

# Summary

- We introduced **J-UniMorph**, a Japanese morphology dataset
  - Based on the UniMorph schema
  - Total: 12,687 instances
  - Covers a wide range of verb inflection forms (average: 118 / word)
    - Compared to Wiktionary Edition's 12 / word
- We developed its interactive **Visualizer** → 
- Future work:
  - Adding unaddressed inflection patterns
  - Expanding coverage to nouns and adjectives

# References

- [Sylak-Glassman'16] John Sylak-Glassman. The composition and use of the universal morphological feature schema (unimorph schema). 2016.
- [Goldman+'23] Omer Goldman, Khuyagbaatar Batsuren, Salam Khalifa, Aryaman Arora, Garrett Nicolai, Reut Tsarfaty, and Ekaterina Vylomova. 2023. SIGMORPHON–UniMorph 2023 shared task 0: Typologically diverse morphological inflection. In Proceedings of the 20th SIGMORPHON workshop on Computational Research in Phonetics, Phonology, and Morphology, pages 117–125, Toronto, Canada. Association for Computational Linguistics.
- [Nitta'07] Nihongo Kijutsu Bunpo Kenkyukai. 2007. *Gendai Nihongo Bunpo 3*. Kuroso Publishers.
- [Nitta'09a] Nihongo Kijutsu Bunpo Kenkyukai. 2009. *Gendai Nihongo Bunpo 2*. Kuroso Publishers.
- [Nitta'09b] Nihongo Kijutsu Bunpo Kenkyukai. 2009. *Gendai Nihongo Bunpo 7*. Kuroso Publishers.
- [Hirabayashi+'88] Yoshisuke Hirabayashi and Yumiko Hama. 1988. *Keigo*. Aratake Publishers.
- [Kato+'89] Yasuhiko Kato and Tsutomu Fukuchi. 1989. *Tense, Aspect, and Mood*. Aratake Publishers.
- [Kamiya'01] Taeko Kamiya. 2001. *The handbook of Japanese verbs*. Kodansha.
- [Takami'11] Ken-ichi Takami. 2011. *Ukemi to Shieki*. Kaitakusha.

# Acknowledgements

- This work was supported by RIKEN Special Postdoctoral Researchers Program and JSPS KAKENHI Grant Numbers JP21K21343, JP22H00524.
- While conducting this research, we received valuable comments from members of the Tohoku NLP Group at Tohoku University, Japan. We deeply appreciate their support and insightful contributions.

# Thank you for listening!

- We introduced **J-UniMorph**, a Japanese morphology dataset
  - Based on the UniMorph schema
  - Total: 12,687 instances
  - Covers a wide range of verb inflection forms (average: 118 / word)
    - Compared to Wiktionary Edition's 12 / word
- We developed its interactive **Visualizer** → 
- Future work:
  - Adding unaddressed inflection patterns
  - Expanding coverage to nouns and adjectives

# Appendix

# | GitHub Repository for J-UniMorph



**J-UniMorph**

<https://github.com/cl-tohoku/J-UniMorph>



jpn repository in UniMorph project

<https://github.com/unimorph/jpn>

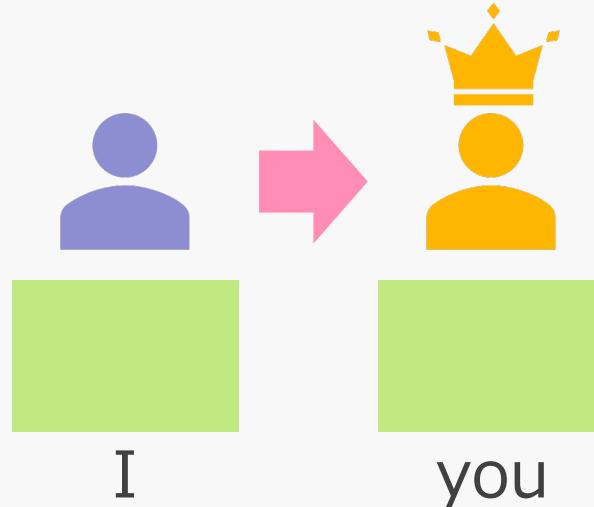
# Introduction: What is UniMorph?

- Universal Morphology (UniMorph) [Sylak-Glassman'16]
  - Collaborative project
  - Covers over **170** languages globally
- UniMorph Feature Schema [Sylak-Glassman'16]
  - Provides standardized morphological features
    - Over 212 feature labels across 23 dimensions of meaning
  - Express with **triplets**: lemma, inflected form, and feature label
    - e.g.

lemma	inflected form	feature label
<i>tabe-ru</i> /食べる	<i>tabe-ta</i> /食べた	V;PST;PFV

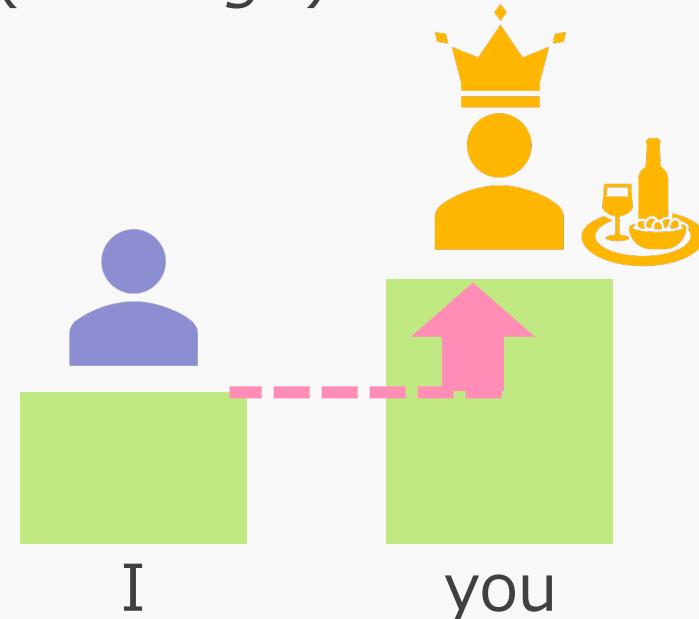
# Another Feature of J-UniMorph: Honorifics

Polite form  
(*Teineigo*)



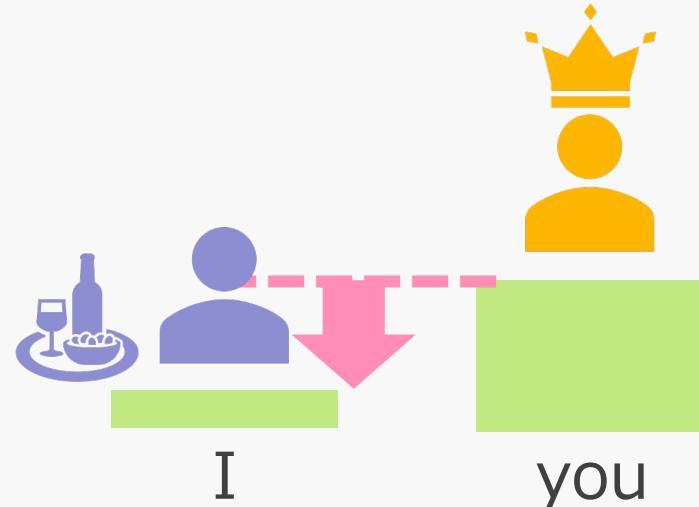
*tabe-ru* eat-V  
*tabe-masu* POL

Respectful form  
(*Sonkeigo*)



*meshiaga-ru* ELEV  
*meshiaga-ri-masu* ELEV;POL

Humble form  
(*Kenjogo*)



*itada-ku* HUMB  
*itada-ki-masu* HUMB;POL

# Honorifics do not correspond one-to-one

- We treated both simple inflections and lexical honorifics
  - Respectful forms are labeled FORM;ELEV (Formal, Referent Elevating)

lemma	inflected form	feature label
<i>tabe-ru</i> (eat)	<i>tabe-ta</i>	V;PST;PFV
<i>tabe-ru</i> (eat)	<i>o-tabे-ni-na-tta</i>	V;FORM;ELEV;PST;PFV
<i>tabe-ru</i> (eat)	<i>meshiaga-tta</i>	V;FORM;ELEV;PST;PFV
<i>tabe-ru</i> (eat)	<i>aga-tta</i>	V;FORM;ELEV;PST;PFV
<i>no-mu</i> (drink)	<i>non-da</i>	V;PST;PFV
<i>no-mu</i> (drink)	<i>meshiaga-tta</i>	V;FORM;ELEV;PST;PFV
<i>no-mu</i> (drink)	<i>aga-tta</i>	V;FORM;ELEV;PST;PFV

# Why did we drop the Derived Verbs?

	Seed Verb (Lemma)	Negative Form
• Regular I	<ul style="list-style-type: none"><li>◦ <i>kak-u</i> (write-V)</li><li>◦ <i>okur-u</i> (send-V)</li></ul>	$\rightarrow$ <i>kak-a-nai</i> $\rightarrow$ <i>okur-a-nai</i>
• Regular II	<ul style="list-style-type: none"><li>◦ <i>tabe-ru</i> (eat-V)</li><li>◦ <i>mi-ru</i> (look-V)</li></ul>	$\rightarrow$ <i>tabe-nai</i> $\rightarrow$ <i>mi-nai</i>
• Irregular	<ul style="list-style-type: none"><li>◦ <i>kuru</i> (come-V)</li><li>◦ <u><i>suru</i></u> (do-V)</li></ul>	$\rightarrow$ <i>konai</i> $\rightarrow$ <u><i>shinai</i></u>
• Derived Verbs	<ul style="list-style-type: none"><li>◦ <i>shiyo-</i><u><i>suru</i></u> (use-N + do)</li><li>◦ <i>benkyo-</i><u><i>suru</i></u> (study-N + do)</li></ul>	$\rightarrow$ <i>shiyo-</i> <u><i>shinai</i></u> $\rightarrow$ <i>benkyo-</i> <u><i>shinai</i></u>

# The Feature of Japanese Morphology

- Morphemes do not correspond one-to-one with meanings

*tabe-nai-desu*

... (I) don't eat. (Colloquial)

*tabe-ru*

eat

*nai*

NEG

*desu*

PRS;IPFV;POL

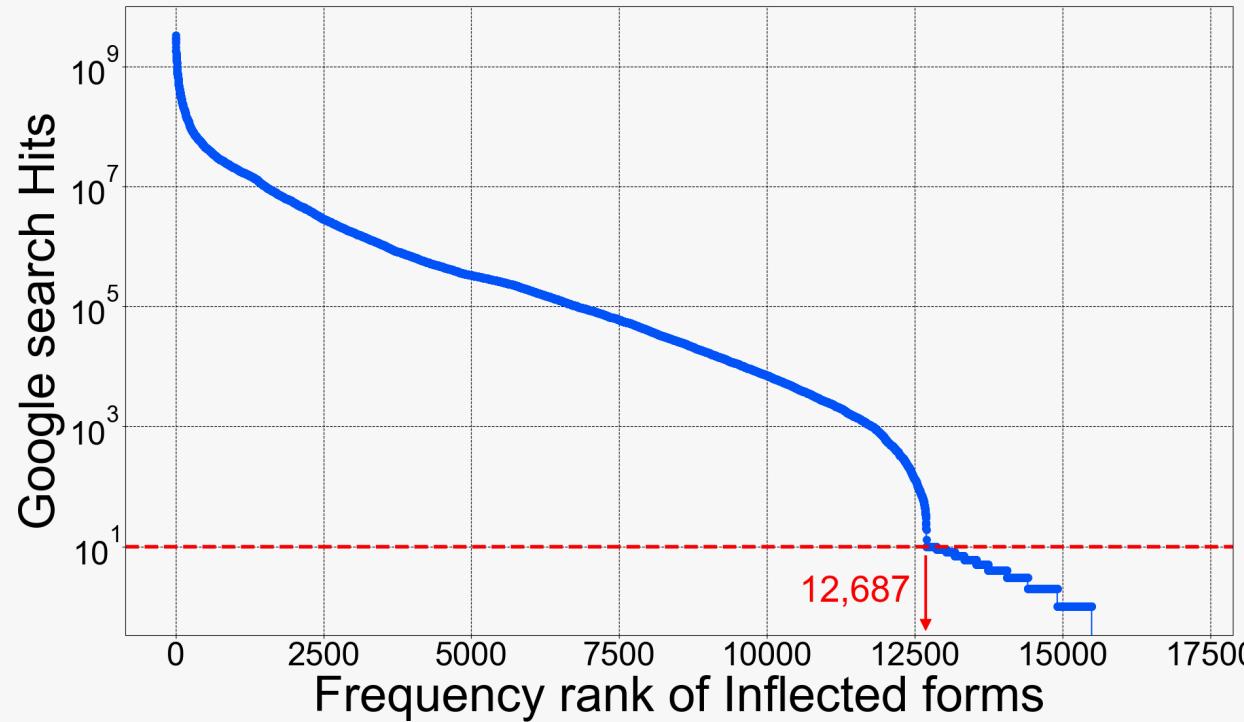
*nai* } *desu*

COL

lemma	inflected form	feature label
<i>tabe-ru</i>	<i>tabe-nai-desu</i>	V;PRS;IPFV;NEG;POL;COL

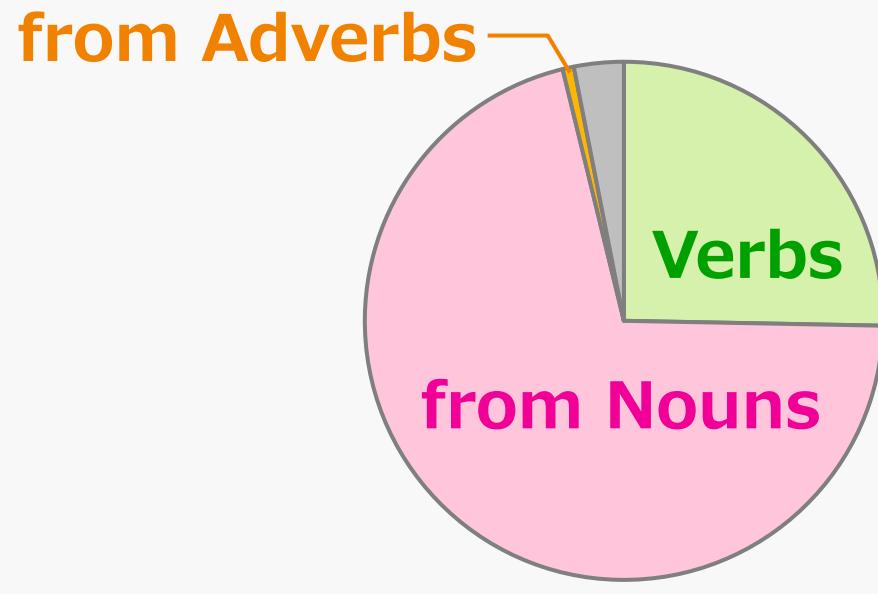
# In the Filtering step: Why 10 hits?

- The trend distinctly shifts when the number of hits reaches 10
- Most of these ( $\leq 10$ ) forms sound unnatural



# Comparison (2/3): Verb Analysis about Derivation

- Wiktionary Edition:
  - About 70% are derived from nouns or adverbs
  - These inflection patterns are morphologically same as “*suru*”



Wiktionary Edition

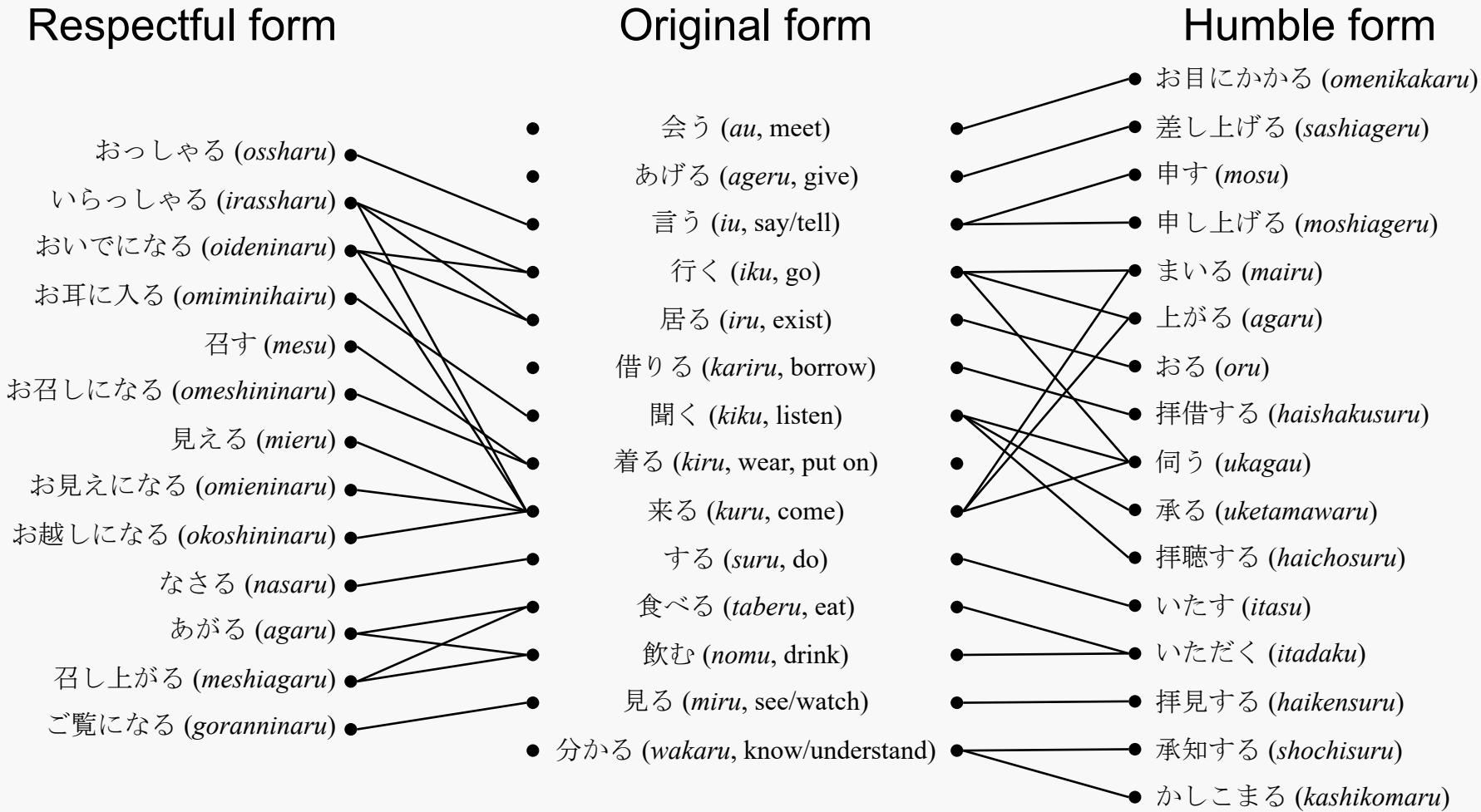


J-UniMorph

Verb from nouns  
are replaced by the  
single verb “*suru*”

# Lexical Honorifics we use

- Morphemes do not correspond one-to-one with meanings



# Randomly retrieve 10 entries (respectively)

Wiktionary Edition

Lemma	Inflected form	Hits
早寝	早寝できる	1,740
鼓する	鼓した	1,610
乗馬	乗馬し	87,600
脱会	脱会し	53,700
鯨飲	鯨飲します	708
断る	断らず	110,000
退部	退部さす	10
管理	管理します	51,600,000
放熱	放熱した	6,890
食指が動く	食指が動こう	369

J-UniMorph (Ours)

Lemma	Inflected form	Hits
話す	お話しください	824,000
渡す	渡させられた	1,220
休む	休まなかつた	88,800
取る	取れます	14,500,000
貸す	貸させていただきます	2,500
走る	走らさない	19,100
乗る	乗らす	19,200
入れる	入れたがりません	219,000
聞く	承りません	9,230,000
撮る	撮らせません	4,720

Hits were assigned by us↑