

Tool untuk menambah Leksem bahasa-bahasa di Indonesia

<https://leksem-indonesia.toolforge.org/>

Penggunaan

1. Memasukkan leksem secara manual

1. Klik Login di <https://leksem-indonesia.toolforge.org/>
2. Pilih salah satu kelas kata, misalnya:
[Leksem kelas nomina](#)
atau
[Leksem kelas nomina \(peng-\)](#)
3. Isi masing-masing bentukan.
4. Apabila ada yang tidak *applicable*, lewati/kosongkan tidak apa-apa.
5. Hasil: <https://www.wikidata.org/wiki/Lexeme:L523837>

1. Memasukkan leksem secara massal

1. Klik Login
2. Pilih salah satu kelas kata, lalu klik di bagian “mode jumlah banyak”
3. Ikuti petunjuk di halaman tersebut. Masing-masing bentukan (*form*) dipisahkan oleh pipa. Apabila ada yang tidak *applicable*, lewati/kosongkan tidak apa-apa. (pipa + pipa)
4. Prasyarat: data harus disiapkan terlebih dahulu (lihat caranya di halaman berikutnya).
5. Pro Tip: mendaftarkan bot dan masuk dengan akun bot sebelum melakukan suntingan massal, atau berpotensi diblokir.

2. Mengapa memasukkan data leksem secara massal?

- Jumlah leksem yang puluhan hingga ratusan ribu per bahasa, tidak logis untuk mengandalkan memasukkan secara manual.
- Lebih cepat, lebih konsisten, lebih terstruktur (tidak ada yang terlewat)
- Lexeme ID-nyaurut dan beraturan 🤖
- Data harus lengkap dan rapi – memakan waktu lama di persiapan data (*data cleanup*)
- Hanya lema saja, tidak perlu definisi. Bisa mengambil dari data kamus mana pun (mis. kamus bahasa daerah)
- Untuk masing-masing bentukan, harus memahami linguistik dan tata bahasa bahasa tersebut terlebih dahulu

3. Mempersiapkan data untuk input massal

1. Siapkan pdf kamus
2. Salin tempel semua isinya, ambil hanya lema (*headword*) saja. Abaikan definisi
3. Menggunakan *spreadsheet* atau program sejenisnya, buat setiap kolom bentukan yang dari lema tersebut
4. Pindahkan ke penyunting teks, ganti tab dengan pipa, simpan sebagai teks (txt)
5. Data siap dimasukkan ke leksem-indonesia.toolforge.org.



3. Memp

1. Siapkan pd
2. Salin temp
3. Mengguna
4. Pindahkan (txt)
5. Data siap d

massal

a. Abaikan

setiap kolom

an sebagai teks

Sejarah

Sejarah singkat mengapa leksem-indonesia lahir dan apa yang membuatnya seperti itu

1. Tahap pertama

April 2021 - Mei 2021

- Dimulai dari usulan penggabungan leksem Indonesia dan Malaysia, yang kontroversial oleh Pengguna:Mahir256 pada bulan April 2021 https://www.wikidata.org/wiki/Wikidata_talk:WikiProject_Indonesia#Lexemes_in_%22Malay%22_and_%22Indonesian%22 (diskusi sangat panjang)
- Kemudian ada keinginan membuktikan perbedaan kedua bahasa, secara leksikal (bukan definisi maupun gramatikal), saya berkesimpulan cara terbaik adalah menunjukkan bahwa ada pengaruh serapan imbuhan dari bahasa gaul, terutama dalam bentuk-bentuk verba (kata kerja), imbuhan-imbuhan nonformal yang tidak dikenal di Melayu.
- Saya menemukan <https://lexeme-forms.toolforge.org/> dan menghubungi Lucas untuk membuat formulir untuk bahasa Indonesia, yang sudah saya persiapkan, https://www.wikidata.org/wiki/Wikidata:Lexeme_Forms/Indonesian namun dihalang-halangi oleh Mahir256 <https://www.wikidata.org/wiki/Topic:W7pbjos433zw8uum>
- Karena tidak ada itikad baik dari keduanya, saya memutuskan untuk melakukan *fork* kode Lucas

2. Tahap kedua

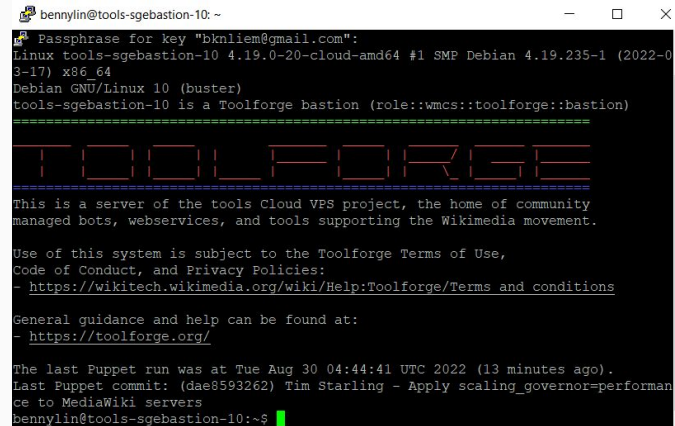
Jun 2021 - Agustus 2021

Supaya mudah, saya putuskan untuk melakukan fork di toolforge. Dokumentasi dapat dibaca di:

- <https://wikitech.wikimedia.org/wiki/Portal:Toolforge>
- <https://wikitech.wikimedia.org/wiki/Help:Toolforge>
- https://wikitech.wikimedia.org/wiki/Portal:Toolforge/Tool_Accounts
- https://wikitech.wikimedia.org/wiki/Help:Create_a_Wiki_media_developer_account#Toolforge_users

Langkah-langkah:

1. meminta akun
2. masuk log dengan PuTTY



```
bennylin@tools-sgebastion-10: ~  
Passphrase for key "bknliem@gmail.com":  
Linux tools-sgebastion-10 4.19.0-20-cloud-amd64 #1 SMP Debian 4.19.235-1 (2022-03-17) x86_64  
Debian GNU/Linux 10 (buster)  
tools-sgebastion-10 is a Toolforge bastion (role::wmcs::toolforge::bastion)  
=====
```

TOOLFORGE

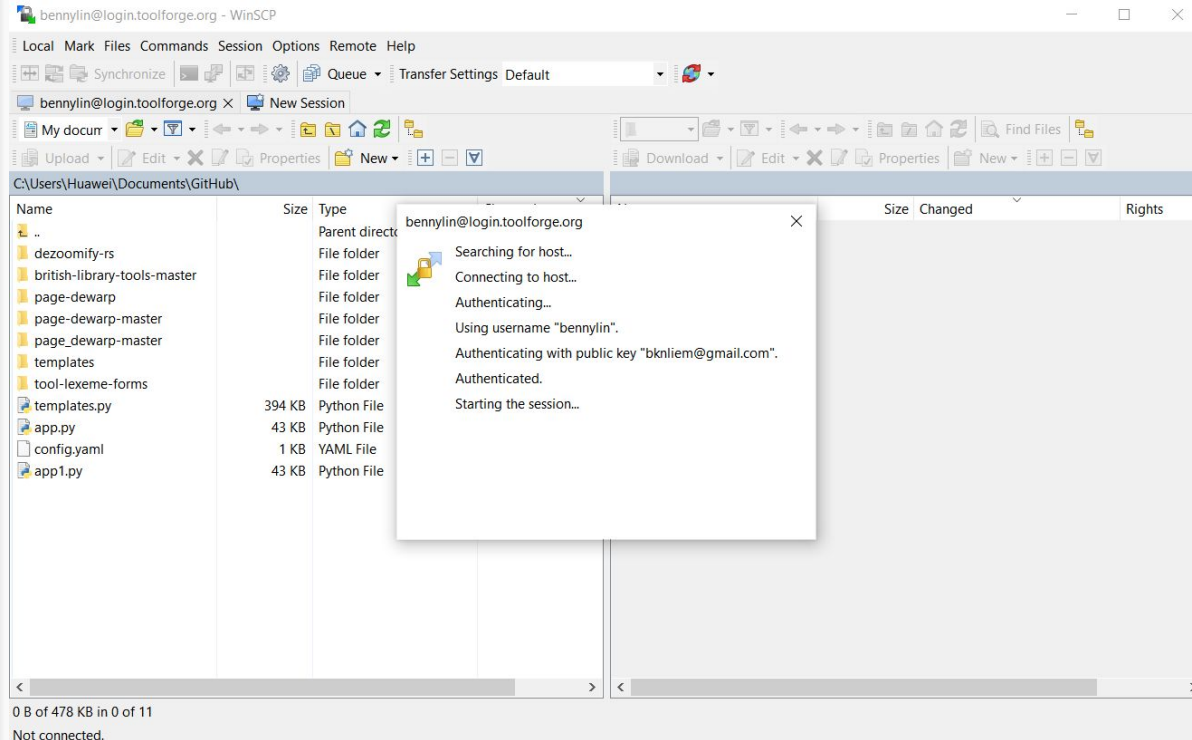
```
-----  
This is a server of the tools Cloud VPS project, the home of community managed bots, webservice, and tools supporting the Wikimedia movement.  
  
Use of this system is subject to the Toolforge Terms of Use, Code of Conduct, and Privacy Policies:  
- https://wikitech.wikimedia.org/wiki/Help:Toolforge/Terms\_and\_conditions  
  
General guidance and help can be found at:  
- https://toolforge.org/  
  
The last Puppet run was at Tue Aug 30 04:44:41 UTC 2022 (13 minutes ago).  
Last Puppet commit: (dae8593262) Tim Starling - Apply scaling_governor=performance to MediaWiki servers  
bennylin@tools-sgebastion-10:~$
```

Tahap kedua (lanj.)

Juni 2021 - Agustus 2021

Langkah-langkah:

1. meminta akun
2. masuk log dengan PuTTY
3. masuk log dengan WinSCP



Tahap kedua (lanj.)

Juni 2021 - Juli 2021

Langkah-langkah:

1. meminta akun
2. masuk log dengan PuTTY
3. masuk log dengan WinSCP, unggah kode lexeme-forms dari Lucas ke
/data/project/leksem-indonesia/www/python/src
4. unggah templates.py dan config.yaml kustom ke root folder
/data/project/leksem-indonesia/
5. di putty, jalankan perintah: "become leksem-indonesia"
lalu "cp templates.py \$HOME/www/python/src/"
Lakukan ini setiap kali ada perubahan templat (ada tambahan bahasa, dll.)
kemudian jalankan perintah "webservice restart"
supaya perubahannya berefek.

3. Tahap ketiga

Juni 2021 - Agustus 2021

(yang nomina tidak selesai, baru sampai di awalan peng-, tidak sempat dilanjutkan karena kehilangan interest / tujuan sudah tercapai, dan tidak sempat menyisihkan waktu lagi)

- Melengkapi bentukan (forms) masing-masing set data
- Melakukan riset
- Membersihkan data
- Memasukkan data lewat “mode jumlah banyak”

Leksem Indonesia (TODO) - nomina dan adjektiva

<https://docs.google.com/spreadsheets/d/1DJNjpbqSwz-z-t0JOpeGvfTqwPptY3hl0d-WCzjDzao/edit#gid=674920041>

Leksem Indonesia (done)-verba, tinggal verba dasar yg belum

https://docs.google.com/spreadsheets/d/1mFu_QGgjpYm4d-D0o1MDe5vYj22zu8q2pU6NDqI0IDk/edit#gid=0

Leksem Jawa, Sunda, Gorontalo, dll.

Leksem Indonesia - non formal:

<https://docs.google.com/spreadsheets/d/1t2aBiJKzXEQyniW1Mj9d2RyNZue-7UjsgKHg5EiEN7k/edit#gid=686325759>

Pembagian dataset

Berdasarkan kelas kata,
dan imbuhan



1. Berdasar kelas kata

Verba: ~15000 leksem

Nomina: ~43000

Adjektiva ~5000

Adverbia ~300

Numeralia ~200, pronomina (~50), partikel (~200): artikula 6, interjeksi 7, interogativa 45, konjungsi 24, preposisi 45,

1. Berdasar kelas kata

Verba: ~15000 leksem

Nomina: ~43000

Adjektiva ~5000

Adverbia ~300

Numeralia ~200, pronomina (~50), partikel (~200): artikula 6, interjeksi 7, interogativa 45, konjungsi 24, preposisi 45,

2. Nomina

Nomina: ~43000

- nomina dasar: 18000 leksem
- nomina majemuk dasar & turunan: 18000 leksem
- -an: ~1600 leksem nomina
- peng-: ~1500 leksem nomina
- peng-an: ~1500 leksem nomina
- per-an: ~500 leksem nomina
- ke-an: ~1400 leksem nomina
- se-: ~300 leksem nomina
- pe-, per-, keber-, keter, kepeng-, kese-, ketidak- (an)

2. Nomina

Nomina: ~43000

- nomina dasar: 18000 leksem
- nomina majemuk dasar & turunan: 18000 leksem
- -an: ~1600 leksem nomina
- **peng-: ~1500 leksem nomina**
- peng-an: ~1500 leksem nomina
- per-an: ~500 leksem nomina
- ke-an: ~1400 leksem nomina
- se-: ~300 leksem nomina
- pe-, per-, keber-, keter, kepeng-, kese-, ketidak- (an)

- Russian (101423)
- Estonian (83208)
- English (71657)
- Malayalam (63315)
- Swedish (36851)
- Latin (32183)
- Hebrew (29912)
- German (27488)
- Basque (22931)
- Spanish (20959)
- Bokmål (17431)
- Slovak (16475)
- Ukrainian (15967)
- Danish (14907)
- Czech (14193)
- French (13707)
- Indonesian (13119)

3. Nomina (peng-)

Nomina: peng-: ~1500 leksem nomina

- baru 100 dari 1282 yang sempat diunggah
- sisanya <https://www.wikidata.org/wiki/User:Bennylin/peng->
- kita kerjakan sekarang

3. Nomina (peng-)

Nomina: peng-: ~1500 leksem nomina

- baru 100 dari 1282 yang sempat diunggah
- sisanya <https://www.wikidata.org/wiki/User:Bennylin/peng->
- kita kerjakan sekarang

Tahap selanjutnya

- https://www.wikidata.org/wiki/User:Bennylin/Leksem_Indonesia

Selanjutnya

1. Membahas topik-topik yang masih didiskusikan

https://www.wikidata.org/wiki/Wikidata:Lexicographical_data/Documentation/Languages/id#Lemma

1. Spelling variants
2. Language codes
3. Regional languages & scripts!
4. Statements
5. Forms
6. Senses, gloss, definitions
7. Etc.: audio, pronunciation, examples, translations

2. Cobalah Query

https://www.wikidata.org/wiki/Wikidata:Lexicographical_data/Documentation/Languages/id#Queries

1. Indonesian nouns
2. Indonesian verbs
3. Indonesian adjectives
4. Get all existing Indonesian lexemes
5. Get the count of lexemes in Indonesian belonging to different lexical categories
6. Ordered by newest to oldest creation time

693001	Q wd:L693001	monumen	Q wd:Q1084	noun
692552	Q wd:L692552	betina	Q wd:Q1084	noun
692346	Q wd:L692346	bahwa	Q wd:Q36484	conjunction
692226	Q wd:L692226	steker	Q wd:Q1084	noun
691676	Q wd:L691676	atau	Q wd:Q36484	conjunction
691675	Q wd:L691675	juga	Q wd:Q380057	adverb
690861	Q wd:L690861	ke	Q wd:Q4833830	preposition
690715	Q wd:L690715	sebagai	Q wd:Q4833830	preposition
689613	Q wd:L689613	dari	Q wd:Q4833830	preposition
689607	Q wd:L689607	di	Q wd:Q4833830	preposition
689600	Q wd:L689600	dan	Q wd:Q36484	conjunction
619093	Q wd:L619093	adaptasi perubahan iklim	Q wd:Q1084	noun
618292	Q wd:L618292	mempergencar	Q wd:Q24905	verb
595062	Q wd:L595062	waktu	Q wd:Q1084	noun
583876	Q wd:L583876	borang	Q wd:Q1084	noun
582682	Q wd:L582682	faffu wasweswos	Q wd:Q34698	adjective
582503	Q wd:L582503	hilir	Q wd:Q1084	noun
582501	Q wd:L582501	hulu	Q wd:Q1084	noun
580939	Q wd:L580939	cepu	Q wd:Q34698	adjective

3. Lihat pula

1. <https://ordia.toolforge.org/language/Q9240> - daftar dan statistik leksem bahasa Indonesia di Wikidata
2. <https://ordia.toolforge.org/language/> - statistik bahasa dengan leksem terbanyak di Wikidata
3. <https://ordia.toolforge.org/search?q=air>
4. <https://machtsinn.toolforge.org/?lang=9240>
5. Recent changes:
<https://www.wikidata.org/wiki/Special:RecentChanges?hidebots=1&hidecategorization=1&tagfilter=OAuth+CID%3A+2271&limit=50&days=30&urlversion=2>

4. Alat-alat menarik

https://www.wikidata.org/wiki/Wikidata:Tools/Lexicographical_data

- <https://orthohin.toolforge.org/add/id> - menambahkan *sense* ke leksem tanpa *sense*
- <https://hangor.toolforge.org/browse/id> - semua leksem, bisa difilter dengan/tanpa *sense*

Hauki is a proof-of-concept tool for searching and displaying Lexemes in a user-friendly way, using lemmas rather than L-numbers as identifiers. <https://www.wikidata.org/wiki/User:Vesihisi/Hauki>

Indonesian

Search lexemes...

Search

fafifu wasweswos

adjective

[\[L582682\]](#)

Forms

fafifuwasweswos

Senses

- Kalimat yang mengandung banyak istilah teknis yang sulit dimengerti oleh orang awam.

[Add a sense!](#)

Examples (not attached to a Sense)

Sekali lagi, mending usahakan tetap jernih dalam diskusinya agar tidak jadi sekadar fafifuwasweswos.

Sudah, tidak usah diladeni. Apalagi yang fafifu wasweswos campur bahasa inggris.

Menurut saya demarkasi story ini terlalu longitudinal elementer. Greymorian dari Goetia saja tafsirnya tidak Mukhalafatu lil hawadisi. Apalagi nubuat hetero-cis-Apocrypha terhollowfikasi ini sangat ndakik-ndakik nggilani.

4. Alat-alat menarik

https://www.wikidata.org/wiki/Wikidata:Tools/Lexicographical_data

- <https://orthohin.toolforge.org/add/id> - menambahkan *sense* ke leksem tanpa *sense*
- <https://hangor.toolforge.org/browse/id> - semua leksem, bisa difilter dengan/tanpa *sense*
- <https://dicare.toolforge.org/lexemes/party.php>
- <https://dicare.toolforge.org/lexemes/challenge.php>
- <https://ordia.toolforge.org/>
- <https://www.wikidata.org/wiki/Module:Lexeme-en>

Snippet code

```

templates = {
    'nomina-indonesia': {
        '@attribution': {'users': ['Bennylin'], 'title':
'User:Bennylin/Leksem_Indonesia'},
        'label': 'Leksem kelas nomina',
        'language_item_id': 'Q9240',
        'language_code': 'id',
        'lexical_category_item_id': 'Q1084',
        'forms': [
            {
                'label': 'nama leksem',
                'example': 'Contoh: [anak]',
                'grammatical_features_item_ids': ['Q106644026'],
            },
            {
                'label': 'tunggal',
                'example': 'Contoh: [anak]',
                'grammatical_features_item_ids': ['Q110786'],
            },
            {
                'label': 'jamak',
                'example': 'Contoh: [anak-anak]',
                'grammatical_features_item_ids': ['Q146786'],
            },
            {
                'label': 'verba',
                'example': 'Contoh: [beranak]',
                'grammatical_features_item_ids': ['Q106614340'],
            },

```

```

        'nomina-indonesia-peng': {
            '@attribution': {'users': ['Bennylin'], 'title':
'User:Bennylin/Leksem_Indonesia'},
            'label': 'Leksem kelas nomina (peng-)',
            'language_item_id': 'Q9240',
            'language_code': 'id',
            'lexical_category_item_id': 'Q1084',
            'forms': [
                {
                    'label': 'kata benda (peng-)',
                    'example': 'Contoh: [pengajar]',
                    'grammatical_features_item_ids': ['Q106644026'],
                },
                {
                    'label': 'kata dasar',
                    'example': 'Contoh: [ajar]',
                    'grammatical_features_item_ids': ['Q111029'],
                },
                {
                    'section_break': True,
                    'label': 'tunggal',
                    'example': 'Contoh: [pengajar]',
                    'grammatical_features_item_ids': ['Q110786'],
                },
                {
                    'label': 'posesif orang pertama',
                    'example': 'Contoh: [pengajarku]',
                    'grammatical_features_item_ids': ['Q71470598'],
                },
            ],

```

```
'adjektiva-indonesia': {
  '@attribution': {'users': ['Bennylin'], 'title':
'User:Bennylin/Leksem_Indonesia'},
  'label': 'Leksem kelas adjektiva',
  'language_item_id': 'Q9240',
  'language_code': 'id',
  'lexical_category_item_id': 'Q34698',
  'forms': [
  {
  'label': 'nama leksem',
  'example': 'Contoh: [berapi-api]',
  'grammatical_features_item_ids': ['Q3482678'],
  },
  {
  'label': 'kata dasar',
  'example': 'Contoh: [api]',
  'grammatical_features_item_ids': ['Q111029'],
  },
  ],
  },
```

```
'adjektiva-indonesia-2': {
  '@attribution': {'users': ['Bennylin'], 'title':
'User:Bennylin/Leksem_Indonesia'},
  'label': 'Leksem kelas adjektiva (2)',
  'language_item_id': 'Q9240',
  'language_code': 'id',
  'lexical_category_item_id': 'Q34698',
  'forms': [
  {
  'label': 'nama leksem',
  'example': 'Contoh: [cantik]',
  'grammatical_features_item_ids': ['Q3482678'],
  },
  {
  'label': 'bentuk superlatif',
  'example': 'Contoh (ter-): [tercantik]',
  'grammatical_features_item_ids': ['Q1817208'],
  },
  {
  'label': 'nomina',
  'example': 'Contoh kata benda: [kecantikan]',
  'grammatical_features_item_ids': ['Q106614337'],
  },
  {
  'label': 'verba',
  'example': 'Contoh kata kerja: [mempercantik]',
  'grammatical_features_item_ids': ['Q106614338'],
  },
  {
  'label': 'tingkat eksefis',
  'example': 'Contoh (ke-an, terlalu ...): [kebanyakan]',
  'grammatical_features_item_ids': ['Q1385613'],
  },
  },
  ],
  },
```

```
'nomina-jawa': {
  '@attribution': {'users': ['Bennylin'], 'title':
'User:Bennylin/Leksem_Jawa'},
  'label': 'Leksem kelas nomina',
  'language_item_id': 'Q33549',
  'language_code': 'jav',
  'lexical_category_item_id': 'Q1084',
  'forms': [
  {
    'label': 'nama leksem',
    'example': 'Contoh: [anak]',
    'grammatical_features_item_ids': ['Q110786'],
  },
  ],
},
```

```
'adjektiva-jawa': {
  '@attribution': {'users': ['Bennylin'], 'title':
'User:Bennylin/Leksem_Jawa'},
  'label': 'Leksem kelas adjektiva',
  'language_item_id': 'Q33549',
  'language_code': 'jav',
  'lexical_category_item_id': 'Q34698',
  'forms': [
  {
    'label': 'nama leksem',
    'example': 'Contoh: [ayu]',
    'grammatical_features_item_ids': ['Q3482678'],
  },
  ],
},
```

```
'verba-jawa': {
  '@attribution': {'users': ['Bennylin'], 'title':
'User:Bennylin/Leksem_Jawa'},
  'label': 'Leksem kelas verba',
  'language_item_id': 'Q33549',
  'language_code': 'jav',
  'lexical_category_item_id': 'Q24905',
  'forms': [
  {
    'label': 'nama leksem',
    'example': 'Contoh: [tangi]',
    'grammatical_features_item_ids': ['Q179230'],
  },
  ],
},
```


Terima kasih!

Presentasi oleh Benny
untuk PPLL WMID
benny@wikimedia.or.id

30-08-2022




Addendum



Yang lupa diutarakan pada pertemuan

- Perlu tool lain (atau tool yang sudah ada perlu ditingkatkan) untuk bisa menambahkan “statement”, seperti akar kata, bandingkan



<https://www.wikidata.org/wiki/Lexeme:L31534>

(L31534) **pembunuh**  **sunting**

id

Language Indonesia
Lexical category nomina

Pernyataan

akar kata	 bunuh	 sunting
	▼ 0 rujukan	+ tambah rujukan
		+ tambah nilai
		+ tambah pernyataan

Addendum

- Perlu tool lain (atau tool yang sudah ada perlu ditingkatkan) untuk bisa menambahkan “statement”, seperti akar kata, bandingkan

<https://www.wikidata.org/wiki/Lexeme:L31534>

dengan

<https://www.wikidata.org/wiki/Lexeme:L693131>

The screenshot shows the Wikidata page for Lexeme L693131. The main entry is for the Indonesian word "pembuntut" (id). It is categorized as a Lexical category nomina in the Indonesian language. Below the main entry, there are sections for "Pernyataan" (statements), "Senses" (with a "+ add Sense" button), and "Forms". Under the "Forms" section, two forms are listed: L693131-F1 for the word "pembuntut" (id) with grammatical features "bentuk tunggal dan jamak", and L693131-F2 for the word "buntut" (id) with grammatical features "akar kata". Each form has a "sunting" (edit) button and a "+ tambah pernyataan" (add statement) button.

Addendum

- Perlu tool lain (atau tool yang sudah ada perlu ditingkatkan) untuk bisa menambahkan “statement”, seperti akar kata, bandingkan
- Big data, datanya akan dibuat di Wiktionary nantinya
- Datanya bersumber dari Wiktionary
<https://id.wiktionary.org/wiki/Kategori:id:Nomina>
<https://id.wiktionary.org/wiki/Kategori:id:Verba>
dst.
- Pentingnya memasukkan pertama kali sudah benar semua, karena kalau sudah ada/sudah dibuat, sulit untuk menyunting.
- Bisa menggunakan Leksem-Indonesia untuk menyunting, mis.:
<https://leksem-indonesia.toolforge.org/template/nomina-indonesia-peng/edit/L31534>
- Kesulitan utama memproses data dan memasukkan lema bahasa daerah adalah kamus-kamus bahasa daerah kebanyakan tidak dilengkapi dengan kelas kata sementara setiap leksem dikategorikan berdasarkan “kategori leksikal” (= kelas kata)
- [BUG]: Lexeme tidak dapat memiliki 2 leksem yang sama dengan kata dasar yang berbeda, misalnya
<https://www.wikidata.org/wiki/Lexeme:L696755>, <https://www.wikidata.org/wiki/Lexeme:L696801>
<https://id.wiktionary.org/wiki/kelebatan>, <https://id.wiktionary.org/wiki/keluaran>
bisa dari akar kata “kelebat”/”keluaran” maupun “lebat”/”luar” (2 kata dengan “sense” yang berbeda)
- Belum ada proses penghapusan leksem di Wikidata(?)