

# Bitext Lexical Data Resources

Bitext Lexical Data Resources are the most comprehensive and consistent set of language data resources in the world, with support for 100+ languages and dialects. This proprietary data has been developed to meet the highest quality standards in the field of computational linguistics. Bitext data is used in production by some of the world's largest and most successful software companies.

## Features

Bitext data sets are rich with comprehensive features. Each language resource has an array of meta data that are relevant to the unique attributes of each specific language, and the data features are consistent across all languages. This comprehensiveness and richness in data provides unlimited flexibility, adaptability and customizability. Inflectional morphology, derivational morphology, use variants and word formation are just a few of the features that are covered by the data. The inflected words in each data set are provided with applicable meta tags and/or information such as:

- Lemma: The canonical form for the inflected word is provided.
- POS: Part of Speech such as noun, verb, adjective, etc. is defined.
- Voice: Verb form is classified as active or passive.
- Tense: Specifies when the action takes place such as past, present, future, etc.
- Aspect: Indicates whether the action is complete, ongoing, habitual, etc.
- Mood: Modality of the verb form is provided: indicative, subjunctive, imperative, etc.
- Person: Verb or pronoun refers to the first, second or third person.
- Number: State of being singular, dual or plural.
- Gender: Noun, verb or adjective forms are provided, masculine, feminine, neuter, etc.
- Case: The function that the noun or adjective plays within a sentence.
- Degree: An adjective is specified as in its positive, comparative or superlative form.
- Definiteness: Specifies whether a noun or adjective refers to a concrete or general concept.

- Polarity: Indicates whether a verb, adjective or noun is in a negative form.
- Contractions: Shortened form of a word or group of words are provided.
- Pronominal Clitics: Clitic pronouns are identified and tagged.
- Formality: Indicates the social status of the speaker in relation to the context.
- Frequency: Relative frequency of the form based on a large general-purpose corpus.
- Named Entities: Pre-defined entities are tagged as person names, places, organization, etc.
- Offensive: Indicates whether the form might be considered offensive in certain contexts.
- Category: Where applicable, frequently used words will be tagged with one of 30 possible categories: animal, body part, brand name, calendar, chemical element, city, clothing, color, computer, country, demonym, family name, female personal name, fruit/vegetable, georegion, greetings, male personal name, measures, organization, pharmaceutical, plant, professions, public figure, relation, seasons, sport, state, transportation, waterway, weather.

## Lexical Data Feature Matrix

LANGUAGE	ISO	TIER	LEMMA	POS	VOICE	TENSE	ASPECT	MOOD	PERSON	NUMBER	GENDER	CASE	DEGREE	DEFINITENESS / STATE	NEGATIVE	CONTRACTIONS	PRONOMINAL CLITICS	FORMALITY	FREQUENCY	NAMED ENTITIES	OFFENSIVE	CATEGORY
Afrikaans	AF	2	x	x		x	x		x	x	x	x	x			x			x	x	x	
Albanian	SQ	2	x	x	x	x		x	x	x	x	x		x					x	x	x	x
Amharic	AM	3	x	x		x			x	x	x	x		x	x			x	x	x	x	
Arabic	AR	3	x	x	x	x		x	x	x	x	x		x			x		x	x	x	
Armenian	HY	3	x	x		x		x	x	x		x	x	x			x		x	x	x	
Assamese	AS	3	x	x		x			x	x	x	x			x			x	x	x	x	
Azeri	AZ	3	x	x	x	x		x	x	x		x			x				x	x	x	
Basque	EU	3	x	x		x		x		x		x							x	x	x	
Belarusian	BE	2	x	x		x	x		x	x	x	x	x						x	x	x	
Bengali	BN	3	x	x		x		x	x	x		x		x	x			x	x	x	x	
Bulgarian	BG	2	x	x		x			x	x	x	x		x					x	x	x	x
Burmese	MY	3	x	x															x	x	x	x
Catalan	CA	1	x	x		x		x	x	x	x					x	x		x	x	x	x
Chinese	ZH	3	x	x															x	x	x	
Croatian	HR	2	x	x	x	x	x		x	x	x	x	x	x		x			x	x	x	x
Czech	CS	2	x	x	x	x		x	x	x	x	x			x			x	x	x	x	x
Danish	DA	1	x	x	x	x			x	x	x	x	x	x					x	x	x	
Dutch	NL	1	x	x		x		x	x	x	x					x			x	x	x	
English	EN	1	x	x		x			x	x	x		x			x			x	x	x	
Esperanto	EO	2	x	x		x				x		x				x			x	x	x	
Estonian	ET	3	x	x	x	x		x	x	x		x	x		x				x	x	x	x

<b>Finnish</b>	FI	3	x	x	x	x		x	x	x		x	x				x	x	x	x	x	x
<b>French</b>	FR	1	x	x		x		x	x	x	x					x	x		x	x	x	
<b>Galician</b>	GL	1	x	x		x		x	x	x	x		x				x		x	x	x	
<b>Georgian</b>	KA	3	x	x		x		x	x	x		x							x	x	x	x
<b>German</b>	DE	1	x	x		x		x	x	x	x	x	x			x			x	x	x	
<b>Greek</b>	EL	2	x	x	x	x	x	x	x	x	x	x	x					x	x	x	x	x
<b>Gujarati</b>	GU	3	x	x		x	x		x	x	x	x							x	x	x	
<b>Hebrew</b>	HE	3	x	x		x			x	x	x				x			x		x	x	
<b>Hindi</b>	HI	3	x	x		x		x	x	x	x	x					x		x	x	x	
<b>Hungarian</b>	HU	3	x	x		x		x	x	x		x	x				x		x	x	x	
<b>Icelandic</b>	IS	2	x	x	x	x		x	x	x	x	x	x						x	x	x	x
<b>Indonesian</b>	ID	3	x	x	x		x				x						x		x	x	x	
<b>Irish Gaelic</b>	GA	2	x	x		x		x	x	x	x	x	x			x			x	x	x	x
<b>Italian</b>	IT	1	x	x		x		x	x	x	x					x	x		x	x	x	
<b>Japanese</b>	JP	3	x	x	x	x												x	x	x	x	
<b>Kannada</b>	KN	3	x	x		x		x	x	x	x	x			x				x	x	x	
<b>Kazakh</b>	KK	3	x	x	x	x	x	x	x	x		x	x		x		x		x	x	x	x
<b>Khmer</b>	KM	3	x	x															x	x	x	x
<b>Korean</b>	KO	2	x	x	x	x		x				x						x	x	x	x	
<b>Kyrgyz</b>	KY	3	x	x	x	x	x	x	x	x		x	x		x		x		x	x	x	
<b>Laos</b>	LO	3	x	x		x													x	x	x	
<b>Latvian</b>	LV	2	x	x	x	x	x	x	x	x	x	x	x	x	x				x	x	x	x
<b>Lithuanian</b>	LT	2	x	x		x		x	x	x	x	x	x	x	x				x	x	x	x
<b>Macedonian</b>	MK	2	x	x		x	x		x	x	x		x	x					x	x	x	
<b>Malay</b>	MS	3	x	x	x		x			x			x				x		x	x	x	x
<b>Malayalam</b>	ML	3	x	x		x		x	x	x	x	x			x			x	x	x	x	
<b>Marathi</b>	MR	3	x	x	x	x	x	x		x	x	x						x	x	x	x	x
<b>Mongolian</b>	MN	3	x	x		x	x	x	x	x		x	x		x				x	x	x	
<b>Nepali</b>	NE	3	x	x	x	x			x	x	x	x	x		x		x	x	x	x	x	
<b>Norwegian Bokmal</b>	NB	1	x	x		x			x	x	x	x	x	x					x	x	x	
<b>Norwegian Nynorsk</b>	NN	1	x	x		x			x	x	x	x	x	x					x	x	x	

Oriya	OR	3	x	x		x	x		x	x		x		x				x	x	x	x	x	
Persian / Farsi	FA	3	x	x		x	x	x	x	x			x	x	x		x		x	x	x		
Polish	PL	2	x	x	x	x	x	x	x	x	x	x							x	x	x	x	
Portuguese	PT	1	x	x		x		x	x	x	x						x		x	x	x		
Punjabi	PA	3	x	x		x	x		x	x	x	x							x	x	x		
Romanian	RO	2	x	x	x	x		x	x	x	x	x		x					x	x	x	x	
Russian	RU	2	x	x		x		x	x	x	x	x	x						x	x	x		
Serbian	SR	2	x	x	x	x	x		x	x	x	x	x	x		x			x	x	x	x	
Sindhi	SD	3	x	x	x	x		x	x	x	x	x						x	x	x	x		
Sinhala	SI	3	x	x		x		x	x		x	x		x				x	x	x	x		
Slovak	SK	2	x	x		x	x		x	x	x	x	x		x				x	x	x	X	
Slovenian	SL	2	x	x		x	x	x	x	x	x	x	x	x					x	x	x	x	
Spanish	ES	1	x	x		x		x	x	x	x						x		x	x	x		
Swahili	SW	3	x	x		x			x	x									x	x	x		
Swedish	SV	1	x	x	x	x		x	x	x	x	x	x	x					x	x	x		
Tagalog / Filipino	TL	3	x	x	x		x	x											x	x	x		
Tamil	TA	3	x	x		x			x	x	x	x					x		x	x	x		
Telugu	TE	3	x	x		x	x		x	x	x	x			x				x	x	x		
Thai	TH	3	x	x														x	x	x	x		
Turkish	TR	3	x	x		x		x	x	x		x			x		x		x	x	x		
Ukrainian	UK	2	x	x		x	x	x	x	x	x	x	x				x		x	x	x		
Urdu	UR	3	x	x		x		x	x	x	x	x					x		x	x	x		
Uzbek	UZ	3	x	x	x	x	x	x	x	x		x	x		x		x		x	x	x		
Vietnamese	VI	3	x	x															x	x	x		
Welsh	CY	3	x	x		x	x	x	x	x	x		x						x	x	x	x	
Zulu	ZU	3	x	x		x		x		x	x	x							x	x	x		
TOTAL LANGUAGES		78		78	78	26	70	25	47	63	69	50	55	34	21	19	11	23	15	78	78	78	25

## Lexical Data Feature Matrix

LANGUAGE VARIANT	ISO	TIER -LEXICAL	LEMMA	POS	VOICE	TENSE	ASPECT	MOOD	PERSON	NUMBER	GENDER	CASE	DEGREE	DEFINITENESS/ STATE	NEGATIVE	CONTRACTIONS	PRONOMINAL CLITICS	FORMALITY	FREQUENCY	NAMED ENTITIES	OFFENSIVE	CATEGORY
Arabic (MSA)	AR	3	x	x	x	x		x	x	x	x	x		x			x		x	x	x	
Arabic (Egypt)	AR	3	x	x	x	x		x	x	x	x	x		x			x		x	x	x	x
Arabic (Gulf)	AR	3	x	x	x	x		x	x	x	x	x		x			x		x	x	x	x
Arabic (Najdi)	AR	3	x	x	x	x		x	x	x	x	x		x			x		x	x	x	x
Chinese (Simplified)	ZH	3	x	x															x	x	x	
Chinese (Traditional)	ZH	3	x	x															x	x	x	
Dutch (Netherlands)	NL	1	x	x		x		x	x	x	x					x			x	x	x	
Dutch (Belgium)	NL	1	x	x		x		x	x	x	x					x			x	x	x	
English (USA)	EN	1	x	x		x			x	x	x		x			x			x	x	x	
English (UK)	EN	1	x	x		x			x	x	x		x			x			x	x	x	
English (India)	EN	1	x	x		x			x	x	x		x			x			x	x	x	
Finnish (Standard)	FI	3	x	x	x	x		x	x	x		x	x				x	x	x	x	x	x
Finnish (Colloquial)	FI	3	x	x	x	x		x	x	x		x	x				x	x	x	x	x	x
French (France)	FR	1	x	x		x		x	x	x	x					x	x		x	x	x	
French (Canada)	FR	1	x	x		x		x	x	x	x					x	x		x	x	x	
French (Switzerland)	FR	1	x	x		x		x	x	x	x					x	x		x	x	x	
German (Germany)	DE	1	x	x		x		x	x	x	x	x	x			x			x	x	x	
German (Switzerland)	DE	1	x	x		x		x	x	x	x	x	x			x			x	x	x	
Italian (Italy)	IT	1	x	x		x		x	x	x	x					x	x		x	x	x	
Italian (Switzerland)	IT	1	x	x		x		x	x	x	x					x	x		x	x	x	
Portuguese (Portugal)	PT	1	x	x		x		x	x	x	x						x		x	x	x	
Portuguese (Brazil)	PT	1	x	x		x		x	x	x	x						x		x	x	x	

Spanish (Spain)	ES	1	x	x		x		x	x	x	x						x		x	x	x	
Spanish (North America)	ES	1	x	x		x		x	x	x	x						x		x	x	x	
Spanish (Central America)	ES	1	x	x		x		x	x	x	x						x		x	x	x	
Spanish (Andes)	ES	1	x	x		x		x	x	x	x						x		x	x	x	
Spanish (Southern Cone)	ES	1	x	x		x		x	x	x	x						x		x	x	x	
<b>TOTAL VARIANTS</b>	<b>27</b>		<b>27</b>	<b>27</b>	<b>6</b>	<b>25</b>	<b>-</b>	<b>22</b>	<b>25</b>	<b>25</b>	<b>23</b>	<b>8</b>	<b>7</b>	<b>4</b>	<b>-</b>	<b>12</b>	<b>18</b>	<b>2</b>	<b>27</b>	<b>27</b>	<b>27</b>	<b>5</b>
<b>TOTAL LANGUAGES AND VARIANTS</b>	<b>105</b>		<b>105</b>	<b>105</b>	<b>32</b>	<b>95</b>	<b>25</b>	<b>69</b>	<b>88</b>	<b>93</b>	<b>73</b>	<b>63</b>	<b>41</b>	<b>25</b>	<b>19</b>	<b>23</b>	<b>41</b>	<b>17</b>	<b>105</b>	<b>105</b>	<b>105</b>	<b>30</b>