

Bitext Linguistic Services

Bitext offers a multilingual platform to analyze & tag text at two levels:

- lexical word-level
- syntactic and semantic sentence-level

The platform has been designed for enterprise use: one single API for any language, platform independent (runs on Linux, Windows, Android, iOS), maximum performance/throughput, minimum footprint and easy scalability. As for deployment, customized solutions are available both for cloud or on-premise.

Lexical Level and Lemmatization

At the lexical level, the main component is the lemmatizer, integrated with tools for decompounding or word segmentation (required by some languages to perform proper lemmatization).

The lemmatizer can be additionally packaged to cover a language analysis full pipeline, from sentence segmentation to full parsing, including tools like spell checking.

Both components of the lemmatizer, data and software, can be distributed integrated or separately. All these tools are available in 78 languages and 27 language variants.

Syntactic Level and Parsing

At the syntactic level, the parser is the main component. The parser analyzes the structure of the sentences in the text and is used for tasks like POS Tagging and Phrase Extraction. Additionally, it is used as the base component for various semantic-level tasks like Entity Extraction, Topic-Level Sentiment Analysis or Categorization. We have developed parsers for 21 languages and are always adding new languages.



Linguistic Services Overview

Lexical Services (No Grammar)

Language identification	Detect the language(s) used in each sentence of a longer input text
	Applicable to all languages
	Example: Oui! I love Paris → "Oui!" – French, "I love Paris" – English
Sentence segmentation	Splits text into sentences, according to language-specific punctuation rules.
	Applicable to all languages.
	Example: Hello! How are you doing? → Hello! How are you doing?
Tokenization	Splits a sentence into words, according to language-specific space and punctuation rules.
	Applicable to most languages (except Chinese, Japanese, Vietnamese, Thai)
	Example: How are you doing? → How are you doing ?
Word segmentation (no-space tokenization)	Split text into words for languages that do not use spaces to separate them.
	Applicable to Chinese, Japanese, Vietnamese, Thai
	Example: 把音量调低一点→ 把 音量 调低 一点
Lemmatization	Return the possible roots for a word form
	Applicable to most languages (except Chinese, Vietnamese, Thai)
	Example: running → run
Decompounding	Split compound words/tokens into its individual component words.
	Applicable to German, Dutch, Norwegian, Swedish, Korean
	Example: Rindfleischetikettierung → Rind Fleisch Etikettierung
Spelling	Check if a word is spelled correctly
	Applicable to all languages
	Example: excelent → incorrect



Syntactic and Semantic Services (Grammar and Meaning)

POS Tagging	Return the parts of speech for each word in a sentence
	Applicable to all languages
	Example: He runs back home → "He" – pronoun, "runs" – verb, "back" – preposition, "home" - noun
Phrase Extraction	Returns the constituents (noun phrases, verb phrases) of a sentence
	Applicable to all languages
	Example: John's sister was performing in the theatre → "John's sister" – NP, "was performing" – VP, "in the theatre" – PP
Parsing	Produce a tree with the hierarchical constituent parts of a sentence (words, phrases, clauses)
	Applicable to all languages
Entity Extraction	Detect proper names (people, places) and other special text (phones, URLs)
	Applicable to all languages
	Example: John lives in New York → "John" – person name, "New York" – place
Topic-Based Sentiment	Returns the sentiment and corresponding topic of opinions in text
Analysis	Applicable to all languages
	Example: I hate my old phone → opinion: "hate" (negative), topic: "my old phone"
Topic Detection	Returns topics of opinions in text
	Applicable to all languages
	Example: I hate my old phone → topic: "my old phone"
Categorization	Returns the categories applicable to a text, based on pre-defined rules
	Applicable to all languages
	Example: John is feeling great. → HAPPINESS
	[RULE: feel + great → HAPPINESS]
	Example: John was weeping like a willow. → SADNESS
	[RULE: weep + like + willow → SADNESS]

Other

Bots & Assistants	Evaluation and training data generation for chatbots/assistants
Dots & Assistants	Evaluation and training data generation for charbots/assistants



Linguistic Services Matrix - Languages

		~	LANGID	SENTENCE SEG	OKENIZATION	WORD SEG	EMMATIZATION	COMPOUNDING	SPELLING	S TAGGING	PHRASE	PARSE	ENTITIES	SENTIMENT	CATEGORIZATION	OPIC	BOTS & ASSISTANTS
LANGUAGE	ISO	TIER	LAI	SEI	T0	WO	Ē	DE	SPI	PO	ЬН	PA	И	SEI	CA.	T0	BO.
Afrikaans	AF	2	x	x	x		x		X								
Albanian	SQ	2	X	Х	Х		Х		Х								
Amharic	AM	3	X	X	X		Х		X								
Arabic	AR	3	Х	Х	Х		Х		Х	Х			Х				
Armenian	HY	3	X	х	х		х		х								
Assamese	AS	3	Х	х	х		х		х								
Azeri	AZ	3	х	x	x		х		x								
Basque	EU	3	х	х	х		х		х								
Belarusian	BE	2	х	х	х		х		х								
Bengali	BN	3	х	х	х		х		х								
Bulgarian	BG	2	х	х	х		х		х								
Burmese	MY	3	х	х	х		х		х								
Catalan	CA	1	х	х	x		х		х	х	x	х		x		x	
Chinese	ZH	3	х	х		х	х		х	х		х	х	х		х	
Croatian	HR	2	х	х	х		х		х			х					
Czech	CS	2	х	х	х		х		х	х		х					
Danish	DA	1	х	х	х		х		х	х		х					х
Dutch	NL	1	X	х	х		х	х	х	X	х	х	х	X	х	х	Х
English	EN	1	х	х	х		х		х	х	х	х	х	х	х	х	х
Esperanto	EO	2	х	х	х		х		х								
Estonian	ET	3	Х	Х	Х		Х		Х								



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



Oriya	OR	3	х	x	x		х		х								
Persian / Farsi	FA	3	Х	х	Х		Х		х	Х			Х				
Polish	PL	2	х	х	х		х		х	х		х					
Portuguese	PT	1	х	х	х		х		х	х	х	х	х	х	х	х	х
Punjabi	PA	3	х	х	х		х		х								
Romanian	RO	2	х	х	Х		X		Х	Х		Х					
Russian	RU	2	х	х	х		х		х	х		х	х				
Serbian	SR	2	х	х	Х		X		Х	Х		Х					
Sindhi	SD	3	х	х	х		х		х								
Sinhala	SI	3	х	х	х		х		х								
Slovak	SK	2	х	х	х		х		х	х		х					
Slovenian	SL	2	х	х	Х		х		х	х		Х					
Spanish	ES	1	х	х	х		х		х	х	х	х	х	х	х	х	х
Swahili	SW	3	Х	х	Х		Х		Х								
Swedish	SV	1	х	х	х		х	х	х	х		х	х				х
Tagalog / Filipino	TL	3	х	х	х		x		х								
Tamil	TA	3	х	х	х		х		х								
Telugu	TE	3	х	х	х		х		х								
Thai	TH	3	x	х		x	х		х								
Turkish	TR	3	х	х	х		х		х	х							
Ukrainian	UK	2	х	х	х		x		х	х		х					
Urdu	UR	3	х	х	х		x		х								
Uzbek	UZ	3	х	х	х		х		х								
Vietnamese	VI	3	х	х		х	X		х	х							
Welsh	CY	3	х	х	х		х		х								
Zulu	ZU	3	х	х	х		х		x								
TOTAL LANGUAGES	78		78	78	74	4	78	7	78	28	8	21	15	9	7	8	9



Linguistic Services Matrix – Language Variants

				(0			Z	D N							N O		
LANGUAGE VARIANT	081	TIER	LANGID	SENTENCE SEG	TOKENIZATION	WORD SEG	LEMMATIZATION	DECOMPOUNDING	SPELLING	POS TAGGING	PHRASE	PARSE	ENTITIES	SENTIMENT	CATEGORIZATION	TOPIC	BOTS & ASSISTANTS
Arabic (MSA)	AR	3	х	х	х		х		х	х			х				
Arabic (Egypt)	AR	3	х	X	x		х		х				x				
Arabic (Gulf)	AR	3	х	х	х		х		х				х				
Arabic (Najdi)	AR	3	х	х	х		х		х				х				
Chinese (Simplified)	ZH	3	х	х		х	х		х	х		Х	х	х			
Chinese (Traditional)	ZH	3	х	х		х	х		х	х		х	х	х			
Dutch (Netherlands)	NL	1	Х	Х	х		Х	х	х	х	Х	Х	Х	х	х	х	х
Dutch (Belgium)	NL	1	х	х	х		х	х	х	х	х	х	х	х	x	х	х
English (USA)	EN	1	х	х	х		х		х	х	х	х	х	х	х	х	х
English (UK)	EN	1	х	х	х		х		х	х	х	х	х	х	x	х	х
English (India)	EN	1	х	х	х		х		х	х	х	х	х	х	х	х	х
Finnish (Standard)	FI	3	х	x	x		х	x	х								
Finnish (Colloquial)	FI	3	х	х	х		х	х	х								
French (France)	FR	1	x	x	х		х		х	x	х	x	х	х	x	x	x
French (Canada)	FR	1	х	х	х		х		х	х	х	х	х	х	х	х	х
French (Switzerland)	FR	1	х	х	х		х		х	х	х	х	х	х	x	x	x
German (Germany)	DE	1	x	х	х		х	х	х	х	х	х	х	х	х	х	х
German (Switzerland)	DE	1	x	x	x		x	x	х	x	х	x	x	х	x	х	х
Italian (Italy)	IT	1	х	х	х		х		х	х	х	х	х	х	х	х	х
Italian (Switzerland)	IT	1	х	x	x		х		х	х	х	х	х	х	x	х	х
Portuguese (Portugal)	PT	1	Х	Х	Х		Х		х	Х	х	Х	Х	Х	х	Х	Х



TOTAL VARIANTS	27		27	27	25	2	27	6	27	22	19	21	25	21	19	10	10
Spanish (Southern Cone)	ES	1	x	х	x		x		x	x	x	х	x	x	x	x	x
Spanish (Andes)	ES	1	x	x	x		x		x	x	x	x	x	x	x	x	x
Spanish (Central America)	ES	1	x	x	x		x		x	х	х	x	x	x	x	х	х
Spanish (North America)	ES	1	x	x	x		x		x	х	x	x	x	x	x	x	x
Spanish (Spain)	ES	1	х	x	x		х		х	х	х	x	x	х	x	x	х
Portuguese (Brazil)	PT	1	х	х	х		х		Х	Х	Х	Х	Х	х	Х	х	Х