



## SUPPORTED LANGUAGES – CV23

## Table of Contents

Detected Languages.....	3
Tokenized Languages.....	5

# Detected Languages

The following table lists exhaustively all the languages detected by CloudView. The Tokenized Languages section lists the languages which support semantic treatment.

Note: The Language Detector document processor used by CloudView for language detection is more accurate for some languages than others.

For example, it detects correctly English, Russian, Greek, Hebrew, and Japanese most of the time.

Note that the more consecutive words in a specific language the document has, the more precision you get.

<b>A</b>	<b>B</b>	<b>C</b>	<b>D</b>	<b>E</b>	<b>F</b>	<b>G</b>	<b>H</b>	<b>I</b>	<b>J</b>	<b>K</b>	<b>L</b>
Abkhazian Afar	Bashkir Basque	Catalan Central Khmer	Danish Dutch	English Esperanto	Faroese Fijian	Galician Georgian	Hausa Hebrew	Icelandic Indonesian	Japanese Java	Kalaallisut Kannada	Lao Latin
Afrikaans Albanian Amharic	Belarusian Bengali Bihari languages	Chinese Corsican Croatian	Dzongkha	Estonian	Finnish French	German Greek Guarani	Hindi Hungarian	Interlingua Interlingue Inuktitut		Kashmiri Kazakh Kinyarwanda	Latvian Lingala Lithuanian
Arabic Armenian Assamese Aymara Azerbaijani	Bislama Breton Bulgarian Burmese	Czech			Gujarati			Inupiaq Irish Italian		Kirghiz Korean Kurdish	
<b>M</b>	<b>N</b>	<b>O</b>	<b>P</b>	<b>Q</b>	<b>R</b>	<b>S</b>	<b>T</b>	<b>U</b>	<b>V</b>	<b>W</b>	<b>X</b>
Macedonian Malagasy	Nauru Nepali	Occitan Oriya	Persian Polish	Quechua	Romanian Romansh	Samoan Sango	Tagalog Tajik	Ukrainian Urdu	Vietnamese Volapük	Welsh Western Frisian	Xhosa
Malay Malayalam	Norwegian	Oromo	Portuguese Punjabi		Rundi Russian	Sanskrit Scottish Gaelic	Tamil Tatar	Uyghur Uzbek		Wolof	Yiddish Yoruba
Maltese Maori			Pushto			Serbian Serbo- Croatian	Telugu Thai				Zhuang Zulu
Marathi						Shona	Tibetan				

Moldavian  
Mongolian

Sindhi              Tigrinya  
Sinhala              Tonga  
Slovak              Tsonga  
Slovenian              Tswana  
Somali              Turkish  
Sotho,  
Southern  
Spanish              Twi  
Sundanese  
Swahili  
Swati  
Swedish

# Tokenized Languages

Cloudview tokenizes many languages by default (known as Standard support), and provides additional semantic features.

To get additional semantic features, you can purchase the Extended Languages add-on.

Use the table below to decide whether to analyze the languages in your corpus with standard or Extended Languages (Basis Tech) tokenizers.

For details on setting up tokenizers, see the CloudView Configuration guide > Configuring Data Processing > Tokenizing Text section.

The CloudView Configuration guide also describes the various semantic features available in CloudView.

Flag	Supported by
<b>Std</b>	Standard tokenizer
<b>Jap</b>	Specific CloudView tokenizer. For more advanced semantic treatment, move to <i>Extended Languages</i> .
<b>Zhi</b>	Specific CloudView tokenizer. For more advanced semantic treatment, move to <i>Extended Languages</i> .
<b>Ext</b>	<i>Extended Languages</i> (Basis Tech) tokenizers. These tokenizers are embedded in CloudView and require the <i>Extended Languages</i> license.
<b>Ext+</b>	<i>Extended Languages</i> (Basis Tech) additional tokenizers (open source add-on). These tokenizers are not embedded in CloudView and require the <i>Extended Languages</i> license. This add-on is available from V6R2016x.R4. Download the add-on at <a href="http://www.basistech.com/language-download/">http://www.basistech.com/language-download/</a> and unzip its content to <INSTALLDIR>/resource/all-arch.
-	not supported

	Tokenization & Normalization	Lemmatization	Stemming	De - agglutination	Sentence Boundary Detection	Part of speech	Noun Phrase Detection	Named Entity Extraction
Afrikaans	<b>Std</b>	-	-	-	-	-	-	-
Albanian	<b>Std/ Ext+</b>	<b>Ext+</b>	-	-	<b>Ext+</b>	<b>Std</b>	-	-
Arabic	<b>Std/ Ext</b>	<b>Ext</b>	-	-	<b>Ext</b>	<b>Ext</b>	<b>Ext</b>	<b>Ext</b>
Basque	<b>Std</b>	-	-	-	-	-	-	-
Bengali	<b>Std</b>	-	-	-	-	-	-	-
Breton	<b>Std</b>	-	-	-	-	-	-	-
Bulgarian	<b>Std/ Ext+</b>	<b>Ext+</b>	-	-	<b>Ext+</b>	<b>Std</b>	-	-
Catalan	<b>Std/ Ext+</b>	<b>Ext+</b>	-	-	<b>Ext+</b>	<b>Std</b>	-	-
Chinese	<b>Zhi/ Ext</b>	-	-	<b>Ext</b>	<b>Ext</b>	<b>Ext</b>	<b>Zhi/ Ext</b>	<b>Ext</b>
Croatian	<b>Std/ Ext+</b>	<b>Ext+</b>	-	-	-	<b>Std</b>	-	-
Czech	<b>Std/ Ext</b>	<b>Ext</b>	-	-	<b>Std/ Ext</b>	<b>Ext</b>	-	<b>Ext</b> (only for dates)
Danish	<b>Std/ Ext</b>	<b>Ext</b>	<b>Std</b>	<b>Ext</b>	<b>Std/ Ext</b>	<b>Std</b>	<b>Std</b>	-
Dutch	<b>Std/ Ext</b>	<b>Ext</b>	<b>Std</b>	<b>Std/ Ext</b>	<b>Std/ Ext</b>	<b>Ext</b>	<b>Std/ Ext</b>	<b>Std/ Ext</b>
English	<b>Std/ Ext</b>	<b>Std/ Ext</b>	<b>Std</b>	-	<b>Std/ Ext</b>	<b>Std/ Ext</b>	<b>Std/ Ext</b>	<b>Std/ Ext</b>
Estonian	<b>Std/ Ext+</b>	<b>Ext+</b>	-	-	<b>Std/ Ext+</b>	<b>Std</b>	-	-
Farsi (Persian)	<b>Ext</b>	-	-	-	<b>Ext</b>	-	-	<b>Ext</b>

Finnish	<b>Std/ Ext</b>	-	-	-	<b>Std/ Ext</b>	<b>Std</b>	-	-
French	<b>Std/ Ext</b>	<b>Std/ Ext</b>	<b>Std</b>	-	<b>Std/ Ext</b>	<b>Std/ Ext</b>	<b>Std/ Ext</b>	<b>Std/ Ext</b>
German	<b>Std/ Ext</b>	<b>Std/ Ext</b>	<b>Std</b>	<b>Std/ Ext</b>	<b>Std/ Ext</b>	<b>Ext</b>	<b>Std/ Ext</b>	<b>Std/ Ext</b>
Greek	<b>Std/ Ext</b>	<b>Ext</b>	-	-	<b>Ext</b>	<b>Ext</b>	-	<b>Ext</b> (except places)
Hebrew	<b>Std/ Ext</b>	-	-	-	<b>Ext</b>	<b>Std</b>	-	<b>Ext</b>
Hindi	<b>Std</b>	-	-	-	-	-	-	-
Hungarian	<b>Std/ Ext</b>	<b>Ext</b>	-	<b>Ext</b>	<b>Std/ Ext</b>	<b>Ext</b>	-	<b>Ext</b> (except places)
Indonesian	<b>Std/ Ext+</b>	<b>Ext+</b>	-	-	<b>Ext+</b>	<b>Std</b>	-	-
Italian	<b>Std/ Ext</b>	<b>Std</b>	<b>Std</b>	-	<b>Std/ Ext</b>	<b>Std/ Ext</b>	<b>Std/ Ext</b>	<b>Std/ Ext</b>
Japanese	<b>Jap/ Ext</b>	<b>Ext</b>	-	<b>Ext</b>	<b>Ext</b>	<b>Ext</b>	<b>Jap/ Ext</b>	<b>Ext</b>
Korean	<b>Ext</b>	<b>Ext</b>	-	<b>Ext</b>	<b>Ext</b>	<b>Ext</b>	-	<b>Ext</b>
Latin	<b>Std</b>	-	-	-	-	-	-	-
Latvian	<b>Std/ Ext+</b>	<b>Ext+</b>	-	-	<b>Ext+</b>	<b>Std</b>	-	-
Malay	<b>Std/ Ext+</b>	<b>Ext+</b>	-	-	<b>Ext+</b>	<b>Std</b>	-	-
Nepali	<b>Std</b>	-	-	-	-	-	-	-
Norwegian	<b>Std/ Ext</b>	<b>Ext</b>	-	<b>Std/ Ext</b>	<b>Std/ Ext</b>	<b>Std</b>	<b>Std</b>	-
Pashto	<b>Ext</b>	-	-	-	<b>Ext</b>	-	-	<b>Ext</b>
Polish	<b>Std/ Ext</b>	<b>Ext</b>	-	-	<b>Std/ Ext</b>	<b>Ext</b>	<b>Std</b>	<b>Ext</b> (except places)

Portuguese	<b>Std/ Ext</b>	<b>Std/ Ext</b>	<b>Std</b>	-	<b>Std/ Ext</b>	<b>Ext</b>	<b>Std/ Ext</b>	<b>Std/ Ext</b>
Romanian	<b>Std/ Ext</b>	<b>Ext</b>	<b>Std</b>	-	<b>Std/ Ext</b>	<b>Std</b>	<b>Std</b>	-
Russian	<b>Std/ Ext</b>	<b>Std/ Ext</b>	<b>Std</b>	-	<b>Std/ Ext</b>	<b>Ext</b>	<b>Std</b>	<b>Ext</b> (except places)
Serbian	<b>Std/ Ext+</b>	<b>Ext+</b>	-	-	-	<b>Std</b>	-	-
Slovak	<b>Std/ Ext+</b>	<b>Ext+</b>	-	-	<b>Std/ Ext+</b>	<b>Std</b>	<b>Std</b>	-
Slovenian	<b>Std/ Ext</b>	<b>Ext+</b>	-	-	<b>Std/ Ext+</b>	<b>Std</b>	<b>Std</b>	-
Spanish	<b>Std/ Ext</b>	<b>Std/ Ext</b>	<b>Std</b>	-	<b>Std/ Ext</b>	<b>Ext</b>	<b>Std/ Ext</b>	<b>Std/ Ext</b>
Swedish	<b>Std/ Ext</b>	<b>Ext</b>	<b>Std</b>	<b>Ext</b>	<b>Std/ Ext</b>	<b>Std</b>	<b>Std</b>	-
Thai	<b>Ext</b>	-	-	-	-	-	-	-
Turkish	<b>Std/ Ext</b>	-	<b>Std</b>	-	<b>Std/ Ext</b>	<b>Std</b>	-	-
Ukrainian	<b>Std/ Ext+</b>	<b>Ext+</b>	-	-	<b>Ext+</b>	<b>Std</b>	-	-
Urdu	<b>Ext</b>	-	-	-	<b>Ext</b>	-	-	<b>Ext</b>
Welsh	<b>Std</b>	-	-	-	-	-	-	-