**MarkLogic**®

# Language Support

MarkLogic® was designed from the start to handle massive amounts of content, and make that content immediately searchable. MarkLogic makes managing multi-lingual data easy, being language-aware so that you get full comprehension and meaning in all of your data.

## Load and Query in Multiple Languages

➜ **Indexing, tokenizing, and stemming** – MarkLogic will recognize the language of content during load, and then will index, tokenize, and stem the content for querying

➜ **Confidence scores** – Language detection is difficult, which is why MarkLogic returns a list of possible languages with confidence scores so you can further assess and determine the most appropriate language to ingest the document

➜ **Tagging elements** – Before loading, you can specify languages in XML documents at the element level by using the xml:lang attribute

➜ **Custom dictionaries** – You can create custom dictionaries that include stemming and tokenization for proper nouns, technical vocabularies, and spelling variations

➜ **Language-aware search** – All searches in MarkLogic Server are language-aware. All queries search the text, tokenize the search terms, and stem in a particular language

➜ **Language encoding** – MarkLogic translates many different document encodings into UTF-8 upon ingestion, though you can also specify the encoding. Supported encodings include ISO-8859-1, ISO-8859-5, ISO-8859-6, cp1251, cp1252, cp1256, Shift JIS, ISO-2022-JP, ISO 2022-KR, GB18030, Big5, and many others

## Advanced Language Support

Advanced language support includes language-specific tokenization, stemming, and collation for custom ordering of results. MarkLogic includes advanced language support for English, plus 15 other languages, which are available as options. If the language you are interested in is not on our current list please contact us.

| English | German | Arabic | Persian (Farsi) | Portuguese |
|---------|--------|--------|-----------------|------------|
| French | Russian | Chinese (Simplified & Traditional) | Dutch | Norwegian (Nynorsk & Bokmål) |
| Italian | Spanish | Korean | Japanese | |

## Basic Language Support

Basic language support enables basic full-text search and includes tokenization using whitespace-delimiters and punctuation. Our basic language support is available for most languages that use Latin, Cyrillic, Arabic, Greek, or Hebrew script, which includes over 200 languages that each have a million or more native speakers.

| | | | | | |
|---|---|---|---|---|---|
| Abaza | Dogrib | Indonesian | Luba-Lulua | Marshallese | Tagalog |
| Abkhazian | Dungan | Ingush | Lule Sami | Pangasinan | Tagbanwa |
| Achinese | Dyula | Interlingua | Luo | Papiamento | Tahitian |
| Adyghe | Eastern Cham | Inupiaq | Lushootseed | Parsi-Dari | Tamashek |
| Afar | Eastern Frisian | Irish | Luxembourgish | Pashto | Tatar |
| Afrikaans | Efik | Javanese | Macedonian | Pohnpeian | Tausug |
| Akan | English-based Creole or Pidgin | Jju | Madurese | Polish | Tetum |
| Albanian | Erzya | Kabardian | Maguindanao | Prussian | Timne |
| Altaic Language | Esperanto | Kabyle | Makasar | Quechua | Tiv |
| Amo | Estonian | Kalaallisut | Malagasy | Reunion Creole French | Tok Pisin |
| Assyrian Neo-Aramaic | Evenki | Kalmyk | Malay | Rhaeto-Romance | Tokelau |
| Asturian | Ewe | Kamba | Maltese | Romanian | Tonga |
| Atsam | Fang | Kanuri | Mansi | Rundi | Tsonga |
| Avaric | Faroese | Kapampangan | Manx | Russia Buriat | Tswana |
| Aymara | Fijian | Kara-Kalpak | Maori | Samaritan Aramaic | Tumbuka |
| Balinese | Filipino | Karachay-Balkar | Mende | Sami Language | Turkish |
| Bambara | Finnish | Kazakh | Minangkabau | Samoan | Turoyo |
| Bashkir | Finno-Ugrian Language | Khakas | Moksha | Sango | Tuvalu |
| Basque | Fon | Khanty | Moldavian | Sasak | Tuvinian |
| Batak | Friulian | Khasi | Mongo | Scots | Twi |
| Batak Toba | Ga | Kikuyu | Morisyen | Scottish Gaelic | Tyap |
| Belarusian | Gagauz | Kimbundu | Mossi | Selkup | Udihe |
| Bemba | Galician | Kinyarwanda | Muslim Tat | Serer | Udmurt |
| Bini | Ganda | Kirghiz | Nanai | Shona | Uighur |
| Bislama | Gilbertesse | Komi-Permyak | Nauru | Shor | Ukrainian |
| Bosnian | Gorontalo | Komi-Zyrian | Navajo | Sicilian | Ulithian |
| Breton | Greek | Koro | Naxi | Sidamo | Umbundu |
| Buginese | Guarani | Koryak | Ndonga | Skolt Sami | Upper Sorbian |
| Bulgarian | Guianese Creole French | Kosraean | Neapolitan | Slovak | Urdu |
| Bushi | Gwich'in | Kpelle | Nenets | Slovenian | Venda |
| Catalan | Haida | Kuanyama | Niuean | Somali | Vietnamese |
| Cebuano | Haitian | Kumyk | Nogai | Songhai | Volapuk |
| Chamorro | Hanunoo | Ladino | North Ndebele | Soninke | Walloon |
| Chechen | Hawaiian | Lahnda | Norwegian Nynorsk | Sundanese | Yakut |
| Chukot | Hebrew | Lak | Nyamwezi | Susu | Yao |
| Chuukese | Hiligaynon | Latin | Nyanja | Swahili | Yapese |
| Chuvash | Hiri Motu | Latvian | Nyankole | Swati | Yiddish |
| Comorian | Hmong | Lezghian | Occitan | Swedish | Yoruba |
| Cornish | Hopi | Limburgish | Norwegian Nynorsk | Sundanese | Yakut |
| Corsican | Hungarian | Lingala | Nyamwezi | Susu | Yao |
| Croatian | Ibibio | Lisu | Nyanja | Swahili | Yapese |
| Czech | Icelandic | Lithuanian | Nyankole | Swati | Yiddish |
| Danish | Igbo | Low German | Occitan | Swedish | Yoruba |
| Dargwa | Iloko | Lower Sorbian | Oromo | Swiss German | Zulu |
| Dogri | Inari Sami | Luba-Katanga | Palauan | Tabassaran | |

## About MarkLogic

For more than a decade, MarkLogic has delivered a powerful, agile, and trusted Enterprise NoSQL database platform that enables organizations to turn all data into valuable and actionable information. For more information, please visit www.marklogic.com.