# Bitext Linguistic Services Overview

**Lexical services (no grammar)**

| | |
|---|---|
| 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 |
| Decompounding | Split compound words/tokens into its individual component words. |
| | Applicable to German, Dutch, Norwegian, Swedish, Korean… |
| | Example: Rindfleischetikettierung → Rind \| Fleisch \| Etikettierung |
| Lemmatization (ambiguous) | Return the possible roots for a word form |
| | Applicable to most languages (except Chinese, Vietnamese, Thai…) |
| | Example: running → run |
| POS Tagging (ambiguous) | Return the possible parts of speech (and optionally other attributes) of a word |
| | Applicable to all languages |
| | Example: run → verb (infinitive), verb (1[st] person singular, present tense), noun (singular) |
| Inflection | Return all forms of a root word |
| | Applicable to most languages (except Chinese, Vietnamese, Thai…) |
| | Example: run → run, runs, ran, running |
| 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 |
| Spell checking | Check if a word is spelled correctly |
| | Applicable to all languages |
| | Example: excelent → incorrect |
| Spell suggestions | Suggest corrections for incorrectly spelled words |
| | Applicable to all languages |
| | Example: excelent → excellent |

**Syntactic services (grammar)**

| | |
|---|---|
| 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 |
| Offensive language detection | Detect offensive or vulgar expressions in text |
| | Applicable to all languages |
| | Example: tell John to f*ck off → "f*ck off" – offensive |
| Anonymization | Remove sensitive or personal information (PII) from text |
| | Applicable to all languages |
| | Example: My name is John and my account number is 1234567 → My name is XXXX and my account number is XXXX. |
| POS-Tagging (disambiguated) | Return the parts of speech for each word in a sentence |
| | Applicable to all languages |
| | Example: John runs back home → "John" – proper noun, "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 |
| Topic-Based Sentiment Analysis | Returns the sentiment and corresponding topic of opinions in text |
| | Applicable to all languages |
| | Example: I hate my old phone → opinion: "hate" (negative), 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 (low level)**

| | |
|---|---|
| Parsing | Produce a tree with the hierarchical constituent parts of a sentence (words, phrases, clauses…) |
| | Applicable to all languages |