

# ItWac

Das "Italian Web-As-Corpus" besteht aus online verfügbaren Texten, die durch *web crawling* gesammelt wurden.

<b>Sprache</b>	Italienisch
<b>Sprachstufe</b>	Standard
<b>Sprachliche Realisierung</b>	schriftlich
<b>Umfang</b>	ca. 1.5–2 Milliarden Tokens
<b>Medium</b>	Texte von Webseiten, die ausschließlich die Domain .it haben
<b>Geographischer Ursprung</b>	Italien
<b>Form der Daten</b>	Online-Texte
<b>Format</b>	Text, XML
<b>Annotation</b>	lemmatisiert, POS-Tags (automatisch annotiert), das Subkorpus von italienischem Wikipedia wurde zusätzlich mit Semantik und Syntax annotiert
<b>Mögliche Suchabfragen</b>	Mit Sketch Engine oder NoSketchEngine können Wortfrequenzen, n-grams, Konkordanzen usw. erstellt werden
<b>Quelle /Herausgeber</b>	Università di Bologna
<b>Nutzungsvoraussetzungen</b>	Anmeldung erforderlich für Sketch Engine
<b>Link</b>	<a href="https://corpora.dipintra.it/">https://corpora.dipintra.it/</a> (NoSketchEngine) oder <a href="https://www.sketchengine.eu/itwac-italian-corpus/">https://www.sketchengine.eu/itwac-italian-corpus/</a>
<b>Zum Zitieren:</b>	Baroni, Marco, Silvia Bernardini, Adriano Ferraresi & Eros Zanchetta. 2009. The WaCky Wide Web: A Collection of Very Large Linguistically Processed Web-Crawled Corpora. <i>Language Resources and Evaluation</i> 43(3). 209–226.