

# itTenTen16

Das Korpus besteht aus online Texten (knapp 5 Mrd. Wörter), die durch *web crawling* gesammelt wurden. Die Texte wurden dann mit automatisierten Tools aufbereitet. Die 'TenTen' Korpora umfassen über 30 Sprachen.

<b>Sprache</b>	Italienisch
<b>Sprachstufe</b>	Standard
<b>Sprachliche Realisierung</b>	schriftlich
<b>Umfang</b>	ca. 5 Mrd. Wörter
<b>Medium</b>	Online-Texte
<b>Geographischer Ursprung</b>	Italien [Online]
<b>Form der Daten</b>	Texte aus dem Internet, die aufbereitet werden, indem unnötige Informationen gelöscht werden (URLs, Duplikate)
<b>Format</b>	online durchsuchbar
<b>Annotation</b>	lemmatiziert, POS-Tags
<b>Mögliche Suchabfragen</b>	SketchEngine hat mehrere Tools, die unterschiedliche Suchabfragen ermöglichen (Konkordanzen, Synonyme, n-grams usw.)
<b>Quelle/Herausgeber</b>	Jakubík, M., Kilgarriff, A., Ková, V., Rychlý, P., & Suchomel, V., Masaryk University / Lexical Computing
<b>Nutzungsvoraussetzungen</b>	Anmeldung über SSO Universität Potsdam
<b>Link</b>	<a href="https://www.sketchengine.eu/ittenten-italian-corpus/">https://www.sketchengine.eu/ittenten-italian-corpus/</a>
<b>Zum Zitieren:</b>	Jakubík, M., Kilgarriff, A., Ková, V., Rychlý, P., & Suchomel, V. 2013. <a href="#">The TenTen corpus family</a> . <i>7th International Corpus Linguistics Conference CL</i> , 125–127.