

# CLARIAH Shared Development Roadmap: Data Model

# Data Model

CLARIAH Services: Data model

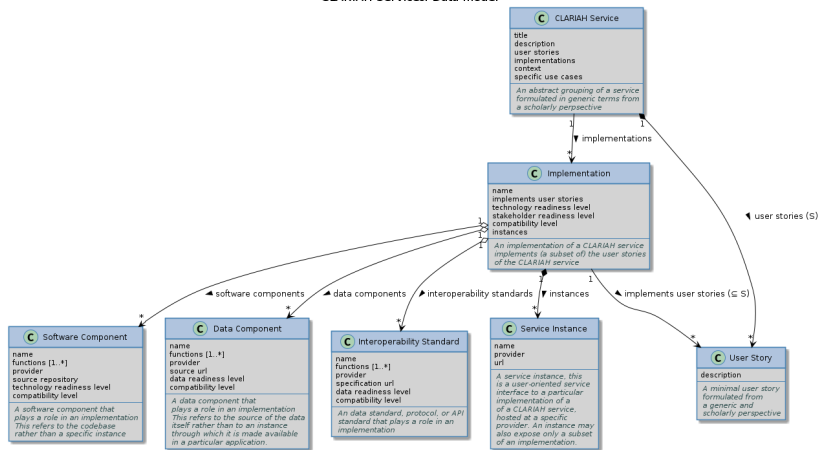


Figure 1: Data model overview

# Template (1/3)

## 2.2.3 Corpus Search: Text & Annotation Search

(Maarten, WP3)

### User story:

**As a scholar**, I want to perform complex searches in text collections/corpora and in the annotations on these collections **in order to** find patterns of specific (often linguistic) constructs for my research purpose.

(2) **As a scholar**, I want to view aggregated results over my results sets, such as distributions, grouped results and statistics **in order to** be able to analyse my data and identify common trends

(3) **As a scholar**, I want to provide my own text collections **in order to** have a platform that enables me to search in them.

(4) **As a scholar**, I want to search in syntactically annotated corpora (treebanks) **in order to** find linguistic patterns for my research purpose. *[this is a more specific instance of the main user story]*

(5) **As a scholar**, I want to automatically enrich my corpus with specific linguistic annotations **in order to** find linguistic patterns for my research purpose.

(6) **As a scholar**, I want uniform and rich access to a large and diverse set of corpora (possibly within a certain domain) **in order to** have a big enough data set to do searches

## Template (2/3)

### Implementations & Software Components

Implementation 1: INT (implements all three stories, might implement 4 in the future. Does not really implement 6)

Component	Function(s)	Instance @Provider
-----------	-------------	-----------------------

Blacklab (using Apache Lucene)	<ul style="list-style-type: none"><li>• Storage engine for text and annotations</li><li>• Query &amp; search engine</li><li>• Indexer to process text corpora with annotations (in specific formats)</li></ul>	<a href="#">AutoSearch @INT</a>  <a href="#">OpenSoNaR @INT</a>
Blacklab Server	<ul style="list-style-type: none"><li>• Web API</li></ul>	
Corpus-fronte nd	<ul style="list-style-type: none"><li>• A search front-end to formulate and execute queries</li><li>• A results front-end to show matches in the corpus, complete with annotations</li><li>• An upload front-end for users to add their own data</li></ul>	
Technology Readiness Level (TRL)	Stakeholder Readiness Level (SRL)	Compatibility Level
8?		

## Template (2/3)

<a href="#">TICCL-tools</a>	Low-level post-OCR normalisation tools that make up the TICCL workflow.	6	UvT
<a href="#">Blacklab</a>	Backend for search over large text collections, including annotations	9	INT
<a href="#">Corpus frontend</a>	Generic search frontend for blacklab	8?	INT
<a href="#">AutoSearch</a>	Specific deployment of Corpus frontend for CLARIAH.	8?	INT
<a href="#">GrETEL</a>	Search in syntactically annotated corpora (treebanks)	8?	UU
<a href="#">PaQu</a>	Search in syntactically annotated corpora (treebanks),	8?	RUG
<a href="#">ELAT</a>	Collaborative web-based linguistic annotation tool (document-based, using FoLiA)	8	KNAW-Huc & CLST RUN (hoster)

Figure 4: Stand-off components