



Text+

Federated Content Search for Lexical Resources (LexFCS) *Specification*

Erik Körner, Thomas Eckart, Uwe Kretschmer, Axel Herold, Frank Wiegand,
Frank Michaelis, Claus Zinn, Matthias Bremm, Louis Cotgrove, Thorsten Trippel,
Felix Rau, Anne Klee, Daniel Werning, Dominik Blöse

Version 0.2, 2024-12-16

Table of Contents

1. Introduction	1
1.1. Terminology	1
1.2. Glossary	1
1.3. Normative References	2
1.4. Non-Normative References	3
1.5. Typographic and XML Namespace conventions	3
2. Summary and Interface Specification	5
2.1. Summary of Changes	5
2.2. Discovery	5
2.2.1. Capabilities	5
2.2.2. Endpoint Description	5
2.3. Searching and Result presentation	5
3. LexFCS Data Model	6
3.1. Entry	6
3.2. Field	6
3.3. Value	8
3.3.1. Generic Value attributes	8
3.3.2. Citation Value attributes	9
3.3.3. Value attributes by Field type	10
3.4. Serialisation	10
4. LexCQL	11
4.1. Search term	11
4.2. Queryable Fields / Indexes	12
4.3. Relations	12
4.4. Relation Modifier	13
4.5. Operators	14
5. LexFCS Data Views	16
5.1. Lexical Data View	16
5.2. Extension of the Hits Data View for LexFCS	16
5.2.1. Using the @kind Attribute	17
A. Normative Appendix	18
A.1. CQL ContextSet specification	18
A.1.1. Indexes	18
A.1.2. Relations	19
A.1.2.1. Implicit Relations	19
A.1.2.2. Defined Relations	19
A.1.3. Booleans	19
A.1.4. Relation Modifiers / Relation Qualifiers	19

A.1.4.1. Functional Modifiers	20
A.1.4.2. Matching	20
A.1.5. Boolean Modifiers	21
A.1.6. Examples	21
B. Non-Normative Appendix	22
B.1. Best Practices	22
B.1.1. Serialization with the Lexical Data View	22
B.1.1.1. Specifying a default language — @xml:lang and @langUri	22
B.1.1.2. Contextualize Value contents for improved interpretability	23
B.1.1.3. Connecting Values using @xml:id and @idRefs	23
B.1.1.4. Connecting Values within Fields using @xml:id and @idRefs to build Value hierarchies	24
B.1.1.5. Referencing external resources using the @ref attribute	25
Changelog	27

Chapter 1. Introduction

The *Lexical Search for Federated Content Search (LexFCS)* specification is an extension of the *CLARIN Federated Content Search (CLARIN-FCS) - Core 2.0* specification that allows search and retrieval of *lexical resources* including dictionaries, encyclopedias, normative data, terminological databases, ontologies etc.

1.1. Terminology

The key words **MUST**, **MUST NOT**, **REQUIRED**, **SHALL**, **SHALL NOT**, **SHOULD**, **SHOULD NOT**, **RECOMMENDED**, **MAY**, and **OPTIONAL** in this document are to be interpreted as in [RFC2119](#).

1.2. Glossary

NOTE | Based on *Glossary* in [FCS Core 2.0](#) specification.

CLARIN-FCS, FCS

CLARIN federated content search, an interface specification to allow searching within resource content of repositories.

Client

A software component, which implements the interface specification to query Endpoints, i.e. an aggregator or a user interface.

CQL

Contextual Query Language, previously known as Common Query Language, is a domain-specific language for representing queries to information retrieval systems, such as search engines, bibliographic catalogs and museum collection databases.

Data View

A Data View is a mechanism to support different representations of search results, e.g. a "hits with highlights" view, an image, or a geolocation.

Endpoint

A software component, which implements the CLARIN-FCS interface specification, and translates between CLARIN-FCS and a search engine.

Hit

Data or a subdivision of data returned by a Search Engine that matches the search criterion. What is considered a Hit highly depends on the Search Engine.

Interface Specification

Common harmonized interface and suite of protocols that repositories are required to implement.

Resource

A searchable and addressable entity offered by an Endpoint, such as a text corpus or a multi-modal corpus.

Result Set

An (ordered) set of Hits that match a search criterion produced by a search engine as the result of processing a query.

Search Engine

A software component within a repository that allows for searching within the repository contents.

SRU

[Search and Retrieve via URL](#) is a protocol for Internet search queries. Originally introduced by Library of Congress [LOC-SRU12](#), later standardization process moved to OASIS [OASIS-SRU12](#), [OASIS-SRU20](#).

1.3. Normative References

NOTE Based on *Normative References* in [FCS Core 2.0](#) specification.

RFC2119

Key words for use in RFCs to Indicate Requirement Levels, IETF RFC 2119, March 1997, <https://www.ietf.org/rfc/rfc2119.html>

XML-Namespaces

Namespaces in XML 1.0 (Third Edition), W3C, 8 December 2009, <http://www.w3.org/TR/2009/REC-xml-names-20091208/>

OASIS-SRU-Overview

searchRetrieve: Part 0. Overview Version 1.0, OASIS, January 2013, <http://docs.oasis-open.org/search-ws/searchRetrieve/v1.0/os/part0-overview/searchRetrieve-v1.0-os-part0-overview.html> (DOC), (PDF)

OASIS-SRU12

searchRetrieve: Part 2. SRU searchRetrieve Operation: APD Binding for SRU 1.2 Version 1.0, OASIS, January 2013, <http://docs.oasis-open.org/search-ws/searchRetrieve/v1.0/os/part2-sru1.2/searchRetrieve-v1.0-os-part2-sru1.2.html> (DOC), (PDF)

OASIS-SRU20

searchRetrieve: Part 3. SRU searchRetrieve Operation: APD Binding for SRU 2.0 Version 1.0, OASIS, January 2013, <http://docs.oasis-open.org/search-ws/searchRetrieve/v1.0/os/part3-sru2.0/searchRetrieve-v1.0-os-part3-sru2.0.html> (DOC), (PDF)

OASIS-CQL

searchRetrieve: Part 5. CQL: The Contextual Query Language version 1.0, OASIS, January 2013, <http://docs.oasis-open.org/search-ws/searchRetrieve/v1.0/os/part5-cql/searchRetrieve-v1.0-os->

[part5-cql.html \(DOC\)](#), [\(PDF\)](#)

LOC-SRU12

SRU Version 1.2: SRU Search/Retrieve Operation, Library of Congress, <http://www.loc.gov/standards/sru/sru-1-2.html>

LOC-CQL

The *Contextual Query Language*, Library of Congress, <https://www.loc.gov/standards/sru/cql/>, see also [OASIS-CQL](#)

LOC-CQLCS

The *CQL Context Set*, Library of Congress, <https://www.loc.gov/standards/sru/cql/contextSets/theCqlContextSet.html>

CLARIN-FCSCore20

CLARIN Federated Content Search (CLARIN-FCS) - Core 2.0 specification, SCCTC FCS Task-Force, <https://office.clarin.eu/v/CE-2017-1046-FCS-Specification-v89.pdf>, ([Sources](#)), ([HTML Build](#)), ([PDF Build](#)),

1.4. Non-Normative References

UD

Universal Dependencies project, that provides terminology for linguistic tags, features and relations, <https://universaldependencies.org/>

EDTF

Extended Date/Time Format (EDTF) Specification, February 4, 2019, The Library of Congress, <https://www.loc.gov/standards/datetime/>

1.5. Typographic and XML Namespace conventions

Sections that are still in discussion and not yet finalized are marked with **(WIP)** and may optionally have some *NOTE* admonition blocks. Details and specifications **SHOULD NOT** be considered stable.

The following typographic conventions for XML fragments will be used throughout this specification:

- `<prefix:Element>`

An XML element with the Generic Identifier *Element* that is bound to an XML namespace denoted by the prefix *prefix*.

- `@attr`

An XML attribute with the name *attr*.

- `string`

The literal *string* **MUST** be used either as the content of an element or value of an attribute.

Endpoints and Clients **MUST** adhere to the [XML-Namespaces](#) specification. The CLARIN-FCS interface specification generally does not dictate whether XML elements should be serialized in their prefixed or non-prefixed syntax, but Endpoints **MUST** ensure that the correct XML namespace is used for elements and that XML namespaces are declared correctly. Clients **MUST** be agnostic regarding syntax for serializing the XML elements, i.e. if the prefixed or un-prefixed variant was used, and **SHOULD** operate solely on *expanded names*, i.e. pairs of *namespace name* and *local name*.

For a list of common XML namespace names and prefixes see the table "XML Namespaces and prefixes" in Section 1.5 of the [FCS Core 2.0 Specification](#).

Chapter 2. Summary and Interface Specification

2.1. Summary of Changes

This specification extends the *CLARIN Federated Content Search (CLARIN-FCS) - Core 2.0* specification in the following ways:

- introducing a new query language, [Chapter 4, LexCQL](#), based on [Contextual Query Language \(CQL\)](#) to allow querying lexical resources,
- extending the basic [Hits Data View](#) for inline markup, and adding a required [Section 5.1, “Lexical Data View”](#) for searching through lexical resources,
- extending the CLARIN-FCS *Endpoint Description* with the *Lexical Search* Capability.

2.2. Discovery

The CLARIN-FCS SRU *explain* response is extended by adding the *Lexical Search* capability and *Lexical Data View* to the *Endpoint Description* to support client auto-configuration.

2.2.1. Capabilities

The *Lexical Search* capability indicates to clients that the FCS endpoint supports searches through lexical resources using [Chapter 4, LexCQL](#) and serializes results in the [Data Views](#).

Table 1. New Lexical Search Capability

Name	Capability Identifier	Summary
<i>Lexical Search</i>	http://clarin.eu/fcs/capability/lex-search	Structured search for lexical resources

2.2.2. Endpoint Description

The *Endpoint Description* remains unchanged. The `<ed:SupportedDataView>` element supports the additional [Section 5.1, “Lexical Data View”](#) which has its own MIME type and `@id` value.

2.3. Searching and Result presentation

Queries **MUST** be formulated using the [Chapter 4, LexCQL](#) query language. Results **MUST** be serialized using both the *Generic Hits* Data View as defined in the CLARIN FCS Core Specification and the [Section 5.1, “Lexical Data View”](#). For the *Generic Hits* Data View, the [Section 5.2, “Extension of the Hits Data View for LexFCS”](#) **MAY** be used, which supports limited inline markup to provide extra context to Hits.

Chapter 3. LexFCS Data Model

The *LexFCS* data model represents information about a single lexical record in a `<lex:Entry>` element and makes no general assumptions about the granularity, type or structure of the underlying lexical resource or its elements.

3.1. Entry

Each `<lex:Entry>` element consists of at least one `<lex:Field>` element, each representing a specific kind of information.

Entry attributes

The `<lex:Entry>` element **MAY** provide the default content language in [IETF BCP 47](#) of the whole entry via the `@xml:lang` attribute. If more information is required, the **OPTIONAL** `@langUri` attribute **MAY** additionally be provided. Language specified on individual `<lex:Value>` elements overrides the `<lex:Entry>` language.

Table 2. Allowed attributes of element `<lex:Entry>`

Value Attribute	Description	Data type	Examples (without quotes)
<code>xml:lang</code>	Content language code in IETF BCP 47	String	"deu", "egy-Egyp", "cop-x-cops"
<code>langUri</code>	URI referencing information about the content language, if not expressable using BCP 47	URI	"https://en.wikipedia.org/wiki/Early_New_High_German"

3.2. Field

Each `<lex:Field>` element has a type that **MUST** be specified using a `@type` attribute. A `<lex:Field>` element of a specific type **MAY** only occur once. The `<lex:Field>` element of type `lemma` is mandatory. For an overview of all allowed field types and their restrictions, refer to [Table 3, “Lexical Search Field types”](#).

The `@type` attribute value of a `<lex:Field>` element governs the semantics, data type, and mandatory / optional attributes of all its `<lex:Value>` elements. A `<lex:Field>` element **MUST** contain at least one `<lex:Value>` element, but **MAY** contain an unlimited number of `<lex:Value>` elements.

Table 3. Lexical Search Field types

Field Type Identifier	Short Description	Value data type	Value examples (without quotes)
<code>lemma</code>	Lemma form, mandatory	String	"Dog", "cat", "walking", "better"
<code>entryId</code>	Identifier of the current entry	String	"104730", "Zwahr-2049-pawk"

Field Type Identifier	Short Description	Value data type	Value examples (without quotes)
phonetic	Phonetic form	String	"tə'dā"
translation	Translation	String	"pechare" (cup)
transcription	Transcription or transliteration	String	-
<i>Prosaic Descriptions</i>			
definition	Definition or description	String	"A car is a vehicle with an engine [...]"
etymology	Etymology information	String	-
<i>Grammar and Morphology</i>			
case	Morphological Case	String	-
number	Morphological Number	String	"sg", "plural"
gender	Morphological Gender	String	"F", "female"
pos	Part-of-Speech	String	"NOUN", "NN", "N", "substantive"
baseform	Baseform (or stem ...) of a lemma	String	
segmentation	Composita segmentation, hyphenation; into phonemes, syllables, subwords	String	"Schach+Brett", "to-day", "ins ta gram men"
sentiment	Sentiment information	String	"angry", "positive", "-0.823"
frequency	Frequency information like occurrences, relative frequency, word rank or frequency class	String	"17234", "0.567"
<i>Relation to other Lexical Entries</i>			
antonym, hyponym, hypernym, meronym, holonym, synonym	Semantic relations	String	-
related	Unspecified relation	String	-
<i>External References</i>			
ref	A URI referencing a related resource.	URI	"http://example.org/somewhere/", "https://coptic-dictionary.org/entry.cgi?ta=C2535"
senseRef	ID of a sense definition.	String	"02961779-n", "4129315-0", "8.10"
<i>Citations / Quotations</i>			

Field Type Identifier	Short Description	Value data type	Value examples (without quotes)
<code>citation</code>	A citation, quotation or example of this entry's lemma.	String	"I got into my car."

3.3. Value

The `<lex:Value>` element contains the actual information content of a specific information type. Its semantics, data type, and allowed / mandatory attributes is governed by the type of the `<lex:Field>` element in which it is contained.

3.3.1. Generic Value attributes

A `<lex:Value>` element can be modified with additional attributes, which are, in most cases, optional but **RECOMMENDED**.

The [Table 4, “Lexical Search Value attributes”](#) lists all allowed attributes for `<lex:Value>` elements. Additional information about their usage and cases in which they are mandatory are stated below.

Table 4. Lexical Search Value attributes

Value Attribute	Description	Data type	Examples (without quotes)
<code>xml:id</code>	XML ID as target for <code>@idRefs</code> attribute	XML ID	"sense1"
<code>xml:lang</code>	Content language code in IETF BCP 47	String	"deu", "egy-Egyp", "cop-x-cops"
<code>langUri</code>	URI referencing information about the content language, if not expressible using BCP 47	URI	"https://en.wikipedia.org/wiki/Early_New_High_German"
<code>preferred</code>	<code><lex:Value></code> is preferred among other <code><lex:Value></code> element among other <code><lex:Value></code> elements in the same <code><lex:Field></code> (may be used for UI hints)	Boolean	"true"
<code>ref</code>	An unspecified reference about the content of the <code><lex:Value></code> element. This can be an audio file for <code>phonetic</code> field value or an external page with more information.	URI	"https://thesaurus-linguae-aegyptiae.de/sentence/ICEDBefDgQaukEMWmajL4HChr4I"
<code>idRefs</code>	Whitespace separated list of XML IDs to reference other field/values	String	"id_sense1 id_sense2"
<code>vocabRef</code>	Reference to the used vocabulary of the content of the <code><lex:Value></code> element	URI	"https://universaldependencies.org/u/pos"

Value Attribute	Description	Data type	Examples (without quotes)
<code>vocabValueRef</code>	Reference to specific value of a vocabulary	URI	"https://universaldependencies.org/u/pos/NOUN", "https://universaldependencies.org/u/feat/Gender.html#Fem"
<code>type</code>	A classification for a field value. It may function to group values. (UI hint)	String	"hyphenation" (segmentation), "sample" (ref)

The attributes `@xml:lang` and `@langUri` are used to specify the object language of the content of the `<lex:Value>` element. They override object language information provided at the `<lex:Entry>` element level, if any. If the `@langUri` attribute is given, the `@xml:lang` attribute is **REQUIRED**.

The `@xml:id` and `@idRefs` attributes are used to link `<lex:Value>` elements inside the same `<lex:Entry>` element, which can be used to highlight these relations in the user interface (like a `citation <lex:Value>` element referencing its corresponding `definition <lex:Value>` element). The semantics of this reference is unspecified.

The `@type` attribute **SHOULD** be used to further specify generic field `<lex:Value>` by using values from limited vocabulary to classify their intended usage. We **RECOMMEND** certain values for different `<lex:Field>` types but it is an open list. Known `@type` values **MAY** be used by Clients to interpret and process `<lex:Value>` contents in specific ways.

The `@vocabRef` attribute refers to the general vocabulary from which the content of the `<lex:Field>` element originates, e.g. <https://universaldependencies.org/u/pos/> for referencing the POS tags of the Universal Dependencies project, which provides context for the interpretation of a string like “noun”. `@vocabValueRef` refers to a specific value of a vocabulary, like <https://universaldependencies.org/u/pos/NOUN> for a noun according to the POS tags of the Universal Dependencies project. If both are provided, they **MUST** refer to the same vocabulary.

3.3.2. Citation Value attributes

Values in `<lex:Field>` elements of type `citation` are allowed more optional attributes, which are listed in Table 5, “Lexical Search Value attributes for type 'citation'”

Table 5. Lexical Search Value attributes for type 'citation'

Value Attribute	Description	Data type	Examples (without quotes)
<code>source</code>	Name of the source from which the citation value was taken	String	"FCS Daily - Latest news and gossip stories"
<code>sourceRef</code>	Reference of the source from which the citation value was taken	URI	"https://www.fcs-daily.de/article/123"
<code>date</code>	Date information for the citation value	EDTF date/time (Level 0)	"1992-02-10", "2024", "2000-01-01T01:02:03Z"

3.3.3. Value attributes by Field type

The Table 6, “Attributes for Lexical Search Field types” lists which attributes **MUST** or **SHOULD** be attached to `<lex:Value>` elements of certain *Field Type Identifiers* to provide Clients context for interpretation of values.

Table 6. Attributes for Lexical Search Field types

Field Type Identifier	Attribute Usage Recommendation
<code>entryId</code>	The attributes <code>@vocabRef</code> or <code>@vocabValueRef</code> MAY be used to provide context.
<code>definition, etymology</code>	Hierarchical relations between <code><lex:Value></code> elements across <code><lex:Field></code> elements can be indicated by using the <code>@xml:id</code> and <code>@idRefs</code> attributes. Subordinate <code><lex:Value></code> elements contain an <code>@idRefs</code> attribute that refers to the <code>@xml:id</code> attribute of another <code><lex:Value></code> element. Clients MAY decide to show this hierarchy.
<code>pos, case, number, gender</code>	The <code>@vocabRef</code> or <code>@vocabValueRef</code> attributes SHOULD be used to allow disambiguation of the provided values. Clients MAY use them to help users, e.g. by providing translations or pointing to definitions.
<code>segmentation</code>	Values SHOULD use the pipe () character as separator to allow for post-processing by Clients (e.g., exchange of separators, splitting of parts). Other separator characters (+, -, ·, ...) MAY be used but uniform processing by Clients can not be guaranteed. The <code>@type</code> attribute SHOULD be used for <code>segmentation <lex:Value></code> elements to indicate what type of segmentation is being performed or described.
<code>ref</code>	The <code>@type</code> attribute MUST be used to indicate what kind of reference is being provided.
<code>senseRef</code>	The attributes <code>@vocabRef</code> or <code>@vocabValueRef</code> MUST be used to provide context.
<code>citation</code>	The <i>special</i> attributes <code>@source</code> , <code>@sourceRef</code> and <code>@date</code> SHOULD be used to provide additional context for <code>citation <lex:Value></code> elements. The <code>@type</code> attribute SHOULD be used to indicate what type of citation, quotation or example is being given.

3.4. Serialisation

The serialisation of a LexFCS entry is specified by the corresponding [Data Views](#).

Chapter 4. LexCQL

LexCQL uses [Contextual Query Language \(CQL\)](#) to query lexical resources available in the FCS. This has the benefit of using an existing, well-known and standardized query language with an established ecosystem, including libraries, parsers and extensive documentation. The proposed Context Set can be found in [Section A.1, "CQL ContextSet specification"](#).

LexCQL queries offer users significant compatibility and flexibility, i.e. a simple query in LexCQL **SHOULD** allow the retrieval of lexical records with alternative spelling or normalisation variants, e.g. upper/lower case, diacritics/umlauts, or other forms of normalization. It should, therefore, be straightforward to formulate meaningful queries, reduce frustration caused by missing or incomplete results, and also enable fuzzy search functionality. Endpoints **SHOULD** support this flexible, user-oriented handling, but are always free to rank more suitable results higher.

However, users **SHOULD** also be given the option of "*sharpening*" search queries using optional operators or modifiers, to refine queries and the associated result sets.

Full support of the LexCQL query language depends on the Endpoint. The Endpoint **MUST** support term-only queries. Endpoints or Clients **MUST** support [Level 2](#) of the CQL server conformance to be able to *parse* (Endpoints) or *serialize* (Clients) all of CQL and respond with appropriate error messages to the search/retrieve protocol interface. For appropriate error messages, refer to [CLARIN-FCS Core 2.0, Section 3.1](#).

NOTE

This does not *imply* that Endpoints are *required* to support all of CQL (and LexCQL), but rather that they are able to *parse* all of CQL and generate the appropriate error message, e.g. if a query includes a feature they do not support.

LexCQL uses the same diagnostics defined in [CLARIN-FCS Core 2.0](#) and [OASIS SRU 2.0](#) to enable endpoints to report error and information messages.

4.1. Search term

A search term **MAY** be enclosed in double quotation marks ("), though it need not be. It **MUST** be enclosed in quotation marks if it contains any of the following characters: left or right angle bracket (<, >), left or right parenthesis ((, (, equal (=), backslash (\), apostrophe/single quotation mark ('), or whitespace. Backslash is used to escape quote and as well as itself.

Examples

- `car`
- `"car"`
- `"car wash"`
- `"car's"`
- `"27\"` to search for term `27"`
- `"\"` to search for term `\`

4.2. Queryable Fields / Indexes

LexCQL allows querying all fields that are supported by the [Chapter 3, LexFCS Data Model](#) and uses their respective [field types](#). Every LexCQL Endpoint **MUST** support queries for field [lemma](#) and **SHOULD** support as many queryable fields as feasible. It is recommended to use the field [lemma](#) as default if no field is specified in the request.

To support querying the language of a lexical record, an additional virtual field [lang](#) has been defined which refers to the language of the whole lexical Entry. It has the same semantics as the [@xml:lang](#) attribute of a [Section 3.2, "Field"](#) in a lexical [Section 3.1, "Entry"](#), where it is queried with the [lang relation modifier](#).

Examples

- [car](#)

Search for lexical records with the term "car" (recommended to be interpreted as query on the field [lemma](#)).

- [lemma = "car"](#)

Search for lexical records with the lemma "car".

- [pos = "NOUN" AND synonym = "house"](#)

Search for nouns that are synonyms of "house".

- [lang = "deu" AND translation =/lang=eng "member of parliament"](#)

Search for an lexical record in German that contains the English translation "member of parliament".

4.3. Relations

LexCQL supports the following relations between field and search term.

- [=](#)

This is the default relation, and the server can choose any appropriate relation or means of comparing the query term with the terms from the data being searched. It is encouraged to interpret this relation in a beginner friendly way and match all reasonably similar terms, based on string normalisation, string similarity, lemmatisation, substring etc.

- [==](#)

This relation is used for exact equality matching. The term in the data is exactly equal to the term in the search query.

- [is](#)

This relation can be used to search for an external value definition specified as URI.

Examples

- `lemma = "car"`

Search for lexical records with the lemma "car". The endpoint decides which records are similar enough and might include in the result set variants, e.g. records with lemma "cars", "CARS" etc.

- `lemma == "car"`

Search for lexical records whose lemma exactly matches the string "car".

- `pos is https://universaldependencies.org/u/pos/NOUN`

Search for lexical records whose part-of-speech is a noun according to the POS tags of the [Universal Dependencies project](https://universaldependencies.org).

4.4. Relation Modifier

Relations *MAY* be modified with relation modifiers, each separated by a slash (/). The following relation modifiers are allowed. The endpoint decides if and to what degree relation modifiers of a query are considered.

- (*masked*, default modifier)

The following patterns and special characters apply for search terms. To explicitly request this functionality, add `q1.masked` as a relation modifier.

- A single asterisk (*) is used to mask zero or more characters.
- A single question mark (?) is used to mask a single character, thus N consecutive question-marks means mask N characters.
- Backslash (\) is used to escape *, ?, quote ("), as well as itself. Backslash not followed immediately by one of these characters is an error.

- `unmasked`

Do not apply masking rules, all characters are literal.

- `lang`

Specifies the language of the requested term. It is encouraged to use and support an [IETF BCP 47](https://tools.ietf.org/html/bcp47) compliant language code.

- `ignoreCase`, `respectCase`

The server is instructed to either ignore or respect the case of the search term, rather than its default behavior (which is unspecified).

- `ignoreAccents`, `respectAccents`

The server is instructed to either ignore or respect diacritics in terms, rather than its default behavior (which is unspecified, but `respectAccents` is recommended).

- `honorWhitespace`

Used with `==` for exact matching to indicate that matching should even include extraneous whitespace (preceding, embedded, or following). In the absence of this modifier it is left to the server to decide whether or not to honor extraneous whitespace.

- `regexp`

The term should be treated as a regular expression. Regular expressions are treated by the individual servers and any features beyond those found in modern POSIX regular expressions will not necessarily be supported by all servers. This modifier overrides the default 'masked' modifier, above.

- `partialMatch`, `fullMatch`

The server is instructed that the search term is either explicitly a partial match (of a potentially longer index value) or should match the index value completely.

WARNING

The CQL Context Set defined in [OASIS](#) specifies the following modifiers that will not be used for LexCQL:
`word`, `string` (how term should be tokenized for matching, `string` for *no break into words*)

Examples

- `lemma = "car s*"`

Search for lexical records whose lemma contains the string "car s" followed by any number of characters (like "car s", "car service" or "car sickness").

- `lemma =/unmasked "car s*"`

Search for lexical records whose lemma is "car s*".

- `synonym =/lang=eng/ignoreCase "handy"`

Search for lexical records with the synonym "handy" (including all case variations, like "HaNdY") in English language.

4.5. Operators

LexCQL supports the following Boolean operators to form complex queries. Boolean operators have the same precedence; they are evaluated left-to-right. Parentheses may be used to override left-to-right evaluation.

- **AND**

The set of records representing two search clauses linked by **AND** is the intersection of the two sets of records representing the two search clauses.

- **OR**

The set of records representing two search clauses linked by **OR** is the union of the two sets of records representing the two search clauses.

- **NOT**

The set of records representing two search clauses linked by **NOT** is the set of records representing the left hand set which are not in the set of records representing the right hand set. **NOT** cannot be used as a unary operator, it is interpreted as "AND NOT".

Examples

- `lemma = "car" AND pos = "NOUN"`

Search for lexical records with both the lemma "car" and the part-of-speech "NOUN".

- `lemma = "car" AND (pos = "NOUN" OR pos = "ADJ")`

Search for lexical records with the lemma "car" and either the part-of-speech "NOUN" or "ADJ".

- `lemma = car NOT pos = "NOUN"`

Search for lexical records with the lemma "car" while not having the part-of-speech "NOUN".

- `pos = NOUN OR verb`

Search for lexical records with the search clauses `pos = NOUN` and `verb`, connected with a Boolean OR operator.

Chapter 5. LexFCS Data Views

Data formats for the representation of results.

5.1. Lexical Data View

The *Lexical (LEX)* Data View is the mandatory serialization of search results for *Lexical Search* queries. It structures information into key and multiple values pairs. More details in [Data Model](#).

Description	The representation of a lexical resource
MIME type	<code>application/x-clarin-fcs-lex+xml</code>
Payload Disposition	<code>inline</code>
Payload Delivery	<code>send-by-default</code> (REQUIRED)
Recommended Short Identifier	<code>lex</code> (RECOMMENDED)
XML Schema	DataView-Lex.xsd

The *Lexical (LEX)* Data View is serialized as XML in the `<fcs:DataView>` element, specified in [FCS Core 2.0 specification](#). (section "Result Format", §2.2.3). The elements `<lex:Entry>`, `<lex:Field>` and `<lex:Value>` with their attributes are the direct serialization of the [Data Model](#).

A minimal example can be seen at [Example of basic Lex Data View](#). A few more targeted examples can be found in [Section B.1.1, "Serialization with the Lexical Data View"](#) demonstrating edge cases and special features when using this Data View.

Example of basic Lex Data View

```
<!-- potential @pid and @ref attributes omitted -->
<fcs:DataView type="application/x-clarin-fcs-lex+xml"
xmlns:fcs="http://clarin.eu/fcs/resource">
  <lex:Entry xmlns:lex="http://textplus.org/fcs/dataview/lex/1.0">
    <lex:Field type="lemma">
      <lex:Value xml:lang="deu">Becher</lex:Value>
    </lex:Field>
    <!-- ... -->
  </lex:Entry>
</fcs:DataView>
```

Serialization examples in the *Lexical (LEX)* Data View can be found at: <https://gitlab.gwdg.de/textplus/ag-fcs-lex-fcs-dataview/-/tree/master/examples>.

5.2. Extension of the Hits Data View for LexFCS

The *Generic Hits (HITS)* Data View is mandatory in [FCS Core 2.0 specification](#) (section "Basic Search", §2.2.3.2). This schema extends the `<hits:Hit>` element with an optional `@kind` attribute (see

Section 5.2.1, “Using the @kind Attribute”), which provides information on the content of the Hit result.

Description	The representation of the hit
MIME type	<code>application/x-clarin-fcs-hits+xml</code>
Payload Disposition	<code>inline</code>
Payload Delivery	<code>send-by-default</code> (REQUIRED)
Recommended Short Identifier	<code>hits</code> (RECOMMENDED)
XML Schema	<code>DataView-LexHits.xsd</code> , based on " <code>DataView-Hits.xsd</code> "

Example of basic Hits Data View

```
<!-- potential @pid and @ref attributes omitted -->
<fcs:DataView type="application/x-clarin-fcs-hits+xml"
xmlns:fcs="http://clarin.eu/fcs/resource">
  <hits:Result xmlns:hits="http://clarin.eu/fcs/dataview/hits">The quick brown
  <hits:Hit>fox</hits:Hit> jumps over the lazy<hits:Hit>dog</hits:Hit>.</hits:Result>
</fcs:DataView>
```

5.2.1. Using the @kind Attribute

To extend the Hits Data View, the `<hits:Hit>` element is reused as per the [FCS Core 2.0 Specification](#), but with the additional optional `@kind` attribute, which provides content hinting. The values of the `@kind` attribute follow the scheme `lex-<FIELD-TYPE>`, e.g. `lex-lemma` for lemma, `lex-pos` for part of speech (POS) tags etc. Textual content outside of `<hits:Hit>` are displayed unchanged.

Example of extended Hits Data View with additional @kind attributes

```
<fcs:DataView type="application/x-textplus-fcs-hits+xml"
xmlns:fcs="http://clarin.eu/fcs/resource">
  <hits:Result xmlns:hits="http://textplus.org/fcs/dataview/hits"><hits:Hit kind="lex-
  lemma">Apple</hits:Hit>: <hits:Hit kind="lex-pos">NOUN</hits:Hit>. <hits:Hit
  kind="lex-def">An apple is an edible fruit produced by an apple
  tree.</hits:Hit></hits:Result>
</fcs:DataView>
```

Endpoints **MUST** generate responses that are valid according to the XML schema "`DataView-LexHits.xsd`".

A. Normative Appendix

A.1. CQL ContextSet specification

Used identifier: <http://text-plus.org/cql/lexres/1.0/> (draft)

Recommended prefix: `lexres`

TIP For more examples of CQL ContextSets, see the [list of Context Sets](#) at the Library of Congress (LoC).

A.1.1. Indexes

TIP For more information about CQL indexes, see [The CQL Context Set, section "INDEXES"](#).

Index title	Description
<code>lemma</code>	Lemma form
<code>entryId</code>	Identifier of the current entry
<code>phonetic</code>	Phonetic form
<code>translation</code>	Translation
<code>transcription</code>	Transcription or transliteration
<code>definition</code>	Definition or description
<code>etymology</code>	Etymology information
<code>case, number, gender</code>	Morphological information like case, number, gender, e.g. in full-form dictionaries
<code>pos</code>	Part of Speech
<code>baseform</code>	Baseform (stem, root, ...) of a lemma
<code>segmentation</code>	Composita segmentation, hyphenation; into phonemes, syllables, subwords
<code>sentiment</code>	Sentiment information
<code>frequency</code>	Frequency information like occurrences, relative frequency, word rank or frequency class
<code>synonym, antonym, hyponym, hypernym, meronym, holonym</code>	Semantic relations, analogous to types in TEI Lex0
<code>related</code>	Unspecified (semantic) relation
<code>ref</code>	A URI referencing a related resource.
<code>senseRef</code>	ID of a sense Sense, Entity, ... definition
<code>citation</code>	A citation, quotation or example of this entry's lemma.

A.1.2. Relations

TIP

More information about CQL relations can be found [The CQL Context Set, section "RELATIONS"](#).

A.1.2.1. Implicit Relations

- =

alias: scr

Different functions based on index. Suggested use:

- *Full match* as default behaviour,
- *Full text search* for longer text fields like **definition**, **etymology** and **citation**.

Endpoints **MAY** decide whether they prefer to focus on precision or recall when matching, e.g. returning case-insensitive matches or other normalizations.

- ==

alias: exact

Exact equality matching. The term in the data is exactly equal to the term in the search query.

A.1.2.2. Defined Relations

- **is**

To query an identifier or entity instead of the string value (=).

- **pos is "https://universaldependencies.org/u/pos/NOUN"** (searching for the <https://universaldependencies.org/u/pos/NOUN> concept)
vs. **pos = "NOUN"** (searching for the **NOUN** string value)

A.1.3. Booleans

TIP

For more information about Booleans in CQL, see [The CQL Context Set, section "BOOLEANS"](#).

- **AND**
- **OR**
- **NOT**
- **PROX** (*is not used!*)

A.1.4. Relation Modifiers / Relation Qualifiers

TIP

For more information about relation modifiers/qualifiers, see [The CQL Context Set](#),

section "RELATION MODIFIERS".

A.1.4.1. Functional Modifiers

- **lang=value**

Specification of a ISO639/BCP47 language tag of the search term. This can for example be used to search for translations in certain languages.

- **ignoreCase, respectCase**

The server is instructed to either ignore or respect the case of the search term, rather than its default behavior (which is unspecified).

- **ignoreAccents, respectAccents**

The server is instructed to either ignore or respect diacritics in terms, rather than its default behavior (which is unspecified, but respectAccents is recommended).

- **partialMatch, fullMatch**

The server is instructed that the search term is either explicitly a partial match (of a potentially longer index value) or should match the index value completely.

A.1.4.2. Matching

- **(masked, default modifier)**

The following patterns and special characters apply for search terms. To explicitly request this functionality, add `q1.masked` as a relation modifier.

- A single asterisk (*) is used to mask zero or more characters.
- A single question mark (?) is used to mask a single character, thus N consecutive question-marks means mask N characters.
- Backslash (\) is used to escape *, ?, quote ("), as well as itself. Backslash not followed immediately by one of these characters is an error.

- **unmasked**

Do not apply masking rules, all characters are literal.

- **honorWhitespace**

Used with `==` for exact matching to indicate that matching should even include extraneous whitespace (preceding, embedded, or following). In the absence of this modifier it is left to the server to decide whether or not to honor extraneous whitespace.

- **regex**

The term should be treated as a regular expression. Regular expressions are treated by the individual servers and any features beyond those found in modern POSIX regular expressions

will not necessarily be supported by all servers. This modifier overrides the default 'masked' modifier, above.

A.1.5. Boolean Modifiers

TIP

More information about Boolean modifiers can be found [The CQL Context Set, section "BOOLEAN MODIFIERS"](#).

none

A.1.6. Examples

1. `cat`
`"cat"`
`"United Nations"`
Different variants to search in the default index.
2. `lemma == mouse`
Search for exact string value `mouse` in `lemma` index.
3. `pos = ADJ`
Search for adjectives.
4. `definition = "cat"`
Search for records whose definition contains the term "cat".
5. `pos = "NOUN" NOT "lion" AND definition = carnivore`
Search for nouns with "carnivore" in definition; exclusion of records with "lion".
6. `pos = NOUN AND (lemma = Apfel OR lemma = "Birne")`
Search for nouns that have lemma with either "Apfel" or "Birne".
7. `translation =/lang=eng car`
Search for lexical entries with an English translation "car".
8. `lemma =/unmasked "^ca?r*"`
Search for the lemma with the literal term "`^ca?r*`".

B. Non-Normative Appendix

B.1. Best Practices

B.1.1. Serialization with the Lexical Data View

B.1.1.1. Specifying a default language — @xml:lang and @langUri

Each `<lex:Value>` element can specify language information using the `@xml:lang` and `@langUri` attributes. To avoid redundancy, a default language **MAY** be specified at the `<lex:Entry>` element. If `<lex:Value>` elements do not specify their language explicitly, they inherit the entry's language information.

`<lex:Value>` elements with the same `@xml:lang` attribute value but different `@langUri` attribute values **MUST** be considered to describe different languages. As a consequence, only `<lex:Value>` elements with the same `@xml:lang` and `@langUri` attribute values as the `<lex:Entry>` element **MAY** inherit the entry's language information. If in doubt, specify the language attributes redundantly.

Default and explicit language specification for Values

```

1 <!-- abbreviated example from examples/dwee-Becher.lex.xml -->
2 <lex:Entry xmlns:lex="http://textplus.org/fcs/dataview/lex/1.0" xml:lang="deu"> ①
3   <lex:Field type="lemma">
4     <lex:Value xml:lang="deu">Becher</lex:Value>
5   </lex:Field>
6   <lex:Field type="segmentation">
7     <lex:Value type="hyphenation">Be|cher</lex:Value> ②
8   </lex:Field>
9   <lex:Field type="hypernym">
10    <lex:Value xml:lang="deu">Gefäß</lex:Value> ③
11    <lex:Value>Gegenstand</lex:Value>
12    <lex:Value xml:lang="goh">pehhari</lex:Value>
13  </lex:Field>
14  <lex:Field type="synonym">
15    <lex:Value>Eimer</lex:Value>
16    <lex:Value xml:lang="deu"
17    langUri="https://en.wikipedia.org/wiki/Early_New_High_German">Kraus</lex:Value> ④
18    <lex:Value xml:lang="goh">kelich</lex:Value>
19  </lex:Field>
20 </lex:Entry>

```

- ① Specifying German as default language of `<lex:Entry>` via `@xml:lang`.
- ② Any `<lex:Value>` element without its own language attributes inherits language information from its parent `<lex:Entry>`, here `@xml:lang="deu"`.
- ③ `@xml:lang="deu"` is redundant.
- ④ Due to `@langUri`, the content language of this `<lex:Value>` element **MUST** be considered to be a different language, even though the `@xml:lang` attribute value is the same as the one at

<lex:Entry> level. As no language information is inherited here, both attributes @xml:lang and @langUri MUST be specified.

B.1.1.2. Contextualize Value contents for improved interpretability

The example [Contextualization of senseRef Values](#) shows how plain text content can be semantically enriched by explicitly stating the used vocabulary or by referencing the value's external definition via attributes @vocabRef or @vocabValueRef.

Contextualization of senseRef Values

```

1 <!-- abbreviated example from examples/wortschatz-Auto.lex.xml -->
2 <lex:Entry xmlns:lex="http://textplus.org/fcs/dataview/lex/1.0" xml:lang="deu">
3   <lex:Field type="lemma">
4     <lex:Value xml:lang="deu">Auto</lex:Value>
5   </lex:Field>
6   <lex:Field type="senseRef">
7     <!-- Dornseiff -->
8     <lex:Value vocabRef="https://doi.org/10.1515/9783110457742">8.10</lex:Value> ①
9     <!-- GermaNet -->
10    <lex:Value vocabRef="http://textplus.sfs.uni-tuebingen.de/api/germanet/synset"
11    >s123456789</lex:Value>
12    <!-- Gemeinsame Normdatei (GND) -->
13    <lex:Value vocabValueRef="http://d-nb.info/gnd/4129315-0">4129315-0</lex:Value>
14    ②
15    <!-- Princeton WordNet -->
16    <lex:Value vocabRef="http://wordnet-rdf.princeton.edu/ontology#Synset"
17    >02961779-n</lex:Value> ③
18  </lex:Field>
19 </lex:Entry>

```

- ① Using the @vocabRef attribute referencing <https://doi.org/10.1515/9783110457742>, the value 8.10 can be interpreted as category "8.10 Auto, Fahrt" of the Dornseiff dictionary, volume 9.
- ② The @vocabValueRef attribute value contains the authoritative reference at the authority file GND for the value 4129315-0 .
- ③ Using the @vocabRef attribute, the value 02961779-n can be interpreted in the context of Princeton WordNet synsets.

B.1.1.3. Connecting Values using @xml:id and @idRefs

Any <lex:Value> element can refer to other <lex:Value> elements in the same <lex:Entry> element by specifying their IDs – stated in their @xml:id attribute – in a whitespace separated list in the @idRefs attribute. The type of this relation is unspecified.

Value relations using @xml:id and @idRefs

```

1 <!-- abbreviated example from examples/GermaNet_Ei.lex.xml -->
2 <lex:Entry xmlns:lex="http://textplus.org/fcs/dataview/lex/1.0" xml:lang="deu">
3   <lex:Field type="lemma">
4     <lex:Value xml:lang="deu">Ei</lex:Value>

```

```

5 </lex:Field>
6 <lex:Field type="senseRef">
7   <lex:Value xml:id="sense_1" vocabRef="http://textplus.sfs.uni-
  tuebingen.de/api/germanet/synset">s39427</lex:Value> ①
8   <lex:Value xml:id="sense_2" vocabRef="http://textplus.sfs.uni-
  tuebingen.de/api/germanet/synset">s25806</lex:Value>
9   <lex:Value xml:id="sense_3" vocabRef="http://textplus.sfs.uni-
  tuebingen.de/api/germanet/synset">s25813</lex:Value>
10 </lex:Field>
11 <lex:Field type="definition">
12   <lex:Value idRefs="sense_1">unzählbar, ohne Plural: [...]</lex:Value> ②
13   <lex:Value idRefs="sense_2">ein Schalengebilde, in dem [...]</lex:Value>
14   <lex:Value idRefs="sense_3">eine Keimzelle</lex:Value>
15 </lex:Field>
16 <lex:Field type="hypernym">
17   <lex:Value idRefs="sense_1">festes Nahrungsmittel</lex:Value> ②
18   <lex:Value idRefs="sense_2">Keim</lex:Value>
19   <lex:Value idRefs="sense_2">Keimling</lex:Value>
20   <lex:Value idRefs="sense_3">Gamet</lex:Value>
21 </lex:Field>
22 <lex:Field type="hyponym">
23   <lex:Value idRefs="sense_2">Brutei</lex:Value>
24   <lex:Value idRefs="sense_2 sense_3">Windei</lex:Value> ③
25 </lex:Field>
26 </lex:Entry>

```

① <lex:Value> element with XML ID `sense_1`.

② <lex:Value> elements referring to the <lex:Value> element with `@xml:id="sense_1"` attribute.

③ A <lex:Value> elements referring to multiple XML IDs.

B.1.1.4. Connecting Values within Fields using `@xml:id` and `@idRefs` to build Value hierarchies

<lex:Value> elements within <lex:Field> elements of type *definition*, *etymology* and *senseRef* can be organized hierarchically by using the `@xml:id` and `@idRefs` attributes. <lex:Value> elements can refer to their head <lex:Value> element via their `@idRefs` attribute.

If `@idRefs` contains multiple IDs then only the first one is assumed to be the direct parent <lex:Value> element, so additional IDs may still connect the <lex:Value> element to other <lex:Value> elements as described in Section B.1.1.3, “Connecting Values using `@xml:id` and `@idRefs`”.

Value hierarchy using `@xml:id` and `@idRefs`

```

1 <!-- abbreviated example from examples/dwdswb-herausgehen.lex.xml -->
2 <lex:Entry xmlns="http://textplus.org/fcs/dataview/lex/1.0" xml:lang="deu">
3   <lex:Field type="lemma">
4     <lex:Value xml:lang="deu">herausgehen</lex:Value>
5   </lex:Field>
6   <lex:Field type="definition">
7     <lex:Value xml:id="def-1">(zu Fuß) (von irgendwo, drinnen) heraus bzw. nach

```

```

draußen, [...]</lex:Value> ①
8   <lex:Value xml:id="def-1-1" idRefs="def-1">etw. (eine Situation, einen
Zustand) hinter [...]</lex:Value> ②
9   <lex:Value xml:id="def-2">sich (aus einem Objekt) herauslösen oder entfernen
lassen</lex:Value>
10  <lex:Value xml:id="def-3">sich ableiten, herausfiltern, schlussfolgern
lassen</lex:Value>
11  <lex:Value xml:id="def-5">besonders von postalischen oder elektronischen
Sendungen; [...]</lex:Value>
12  <lex:Value xml:id="def-5-1" idRefs="def-5">von Informationen; gestreut ,
verbreitet [...]</lex:Value>
13  </lex:Field>
14  <lex:Field type="citation">
15    <lex:Value idRefs="def-1">ich sah sie aus dem Garten herausgehen</lex:Value>
16    <lex:Value idRefs="def-1" source="Die Welt, 29.02.2020">Nach etwa 350 Metern
[...]</lex:Value> ③
17    <lex:Value idRefs="def-1-1" source="Neue Westfälische, 08.05.2023">Die Kliniken
[...]</lex:Value>
18  </lex:Field>
19 </lex:Entry>

```

- ① A `<lex:Value>` element with an `@xml:id` attribute of value `def-1`.
- ② A `<lex:Value>` element stating in its `@idRefs` attribute to be a child of the `<lex:Value>` element with `@xml:id` attribute of value `def-1`.
- ③ `<lex:Value>` elements in `<lex:Field>` elements of type other than `definition` or `etymology` can refer to the same `<lex:Value>` element but will not be seen as children of this hierarchy, only as related in an unspecified way.

B.1.1.5. Referencing external resources using the `@ref` attribute

The `@ref` attribute is used to reference an external resource that provides additional information about the `<lex:Value>` element's content. Depending on the `<lex:Field>` type, the referenced resource may be interpreted differently.

For example in `<lex:Value>` elements in `<lex:Field>` element with `@type="phonetic"`, the external resource might provide a link to an audio file with pronunciation.

Audio integration for phonetic Values

```

1 <!-- abbreviated example from examples/dwdswb-herausgehen.lex.xml -->
2 <lex:Entry xmlns="http://textplus.org/fcs/dataview/lex/1.0" xml:lang="deu">
3   <lex:Field type="lemma">
4     <lex:Value xml:lang="deu">herausgehen</lex:Value>
5   </lex:Field>
6   <lex:Field type="phonetic">
7     <lex:Value ref="https://www.dwds.de/audio/032/herausgehen.mp3">hɛ'raʊsge:ə
n</lex:Value>
8   </lex:Field>

```

9 `</lex:Entry>`

Changelog

2024-12-16 — Proposal of LexFCS (v0.2)

- Specify *Lex Data Model* with *Lex Data View* serialization
- LexCQL for new *Lex Data Model*

2024-04-10 — Publication on CLARIN-ERIC GitHub

2023-05-09 — Publication of first draft (v0.1) on Zenodo

- Propose *LexCQL* query language and extension of *Hits Data View* for LexFCS
- LexFCS CQL Context Set for *LexCQL*