eLex
2023

# A Federated Search and Retrieval Platform for Lexical Resources in Text+ and CLARIN

Thomas Eckart, Axel Herold, **Erik Körner**, Frank Wiegand

Saxon Academy of Sciences and Humanities in Leipzig, Leipzig, Germany
Berlin-Brandenburg Academy of Sciences and Humanities, Berlin, Germany
{eckart,koerner}@saw-leipzig.de, {herold,wiegand}@bbaw.de

NFDI Konsortium

NFDI DIREKTORAT

Sammlungen - Lexikalische Ressourcen - Editionen

Text+

Infrastruktur / Betrieb

NFDI Konsortium

berlin-brandenburgische
AKADEMIE DER WISSENSCHAFTEN

Sächsische Akademie der Wissenschaften
zu Leipzig

# Outline

» Text+ and Lexical Resources

» Federated Content Search Infrastructure

» FCS Specification Extension for Lexical Resources

» Next Steps and Future Work

A Federated Search and Retrieval Platform for Lexical Resources in Text+ and CLARIN
eLex 2023: electronic lexicography in the 21st century

2

# Text+ and NFDI

» [Text+](): research data consortium focused on **language and text data**

» part of Germany's National Research Data Infrastructure ([NFDI]())

    » Aims: make research data **available** for scientific usage, support their **interlinkage**, and their long-term **preservation**

    » inter-disciplinary network of data and services based on common standards and the FAIR principles between consortia from various research areas

A Federated Search and Retrieval Platform for Lexical Resources in Text+ and CLARIN
eLex 2023: electronic lexicography in the 21st century

3

# Text+ Data Domains

» 3 data domains: Collections, Lexical Resources, Editions

» Joint working groups



A Federated Search and Retrieval Platform for Lexical Resources in Text+ and CLARIN
eLex 2023: electronic lexicography in the 21st century

# Data Domain Lexical Resources

» Thematic clusters

- » German dictionaries in the European context
- » Born-digital lexical resources
- » Non-latin scripts and under-resourced languages

» Resource types

- » Dictionaries (mono/bilinigual)
- » Encyclopedias
- » Normative Data (GND, …)
- » Terminology Databases
- » Ontologies, Knowledge Databases
- » Word Nets (Princeton, GermaNet)
- » …

A Federated Search and Retrieval Platform for Lexical Resources in Text+ and CLARIN
eLex 2023: electronic lexicography in the 21st century

5

# Data Domain Lexical Resources
# Data Formats

» Wide, diverse spread of formats with custom search functionalities

   » Generic and customized TEI/XML to legacy XML formats

   » Table-like serializations (lemma lists, frequency information)

   » Custom, proprietary formats

   » Geographic information (images of maps), Character Sets

   » ...

» Challenges due to heterogeneity for unified representation for search and retrieval

A Federated Search and Retrieval Platform for Lexical Resources in Text+ and CLARIN
eLex 2023: electronic lexicography in the 21st century

6

# Data Domain Lexical Resources Findability

» **Decentralized** dictionary platform, **federated** approach

   » Heterogeneous nature of resources, formats, annotation levels, technical architectures

» CLARIN *Federated Content Search (FCS)*

   » Framework for accessing spatially distributed text corpora

   » Common specification of techn. interfaces, data formats, query languages

A Federated Search and Retrieval Platform for Lexical Resources in Text+ and CLARIN
eLex 2023: electronic lexicography in the 21st century

7

# Other Approaches for Linking Lexical Resources

» Often organisationally restricted, e.g. "Wörterbuchnetz" by Trier, Global WordNet Association

» Initiatives in common research infrastructures (CLARIN, DARIAH); ELEXIS

» Standardized formats, e.g. TEI, refinements by DARIAH, ELEXIS; RDF/OntoLex, …

» Collaborative / not exclusively academic dictionaries, e.g. Wiktionary, DBPedia, Wikidata

A Federated Search and Retrieval Platform for Lexical Resources in Text+ and CLARIN
eLex 2023: electronic lexicography in the 21st century

8

# Why Federated Content Search?

» Advantages to endpoints (data centers)

  » Full control about access to resources

    » Copyright, licensing or data protection

  » Knowledge about resource, e.g. how to search, rank results, ...

  » Visibility

» Advantages to end users

  » Ease of use and simple overview of resources and results

  » Simple search box – „Just like google"

  » For details and expert search options → backlink to endpoint

A Federated Search and Retrieval Platform for Lexical Resources in Text+ and CLARIN
eLex 2023: electronic lexicography in the 21st century

9

# CLARIN Federated Content Search (FCS)

» Federated "Corpus Query Platform"

» [FCS](#) =



    » RESTful protocol

    » Query languages & data formats

    » Data Aggregator + web portal

    » (Software ecosystem)

» FCS ≠ complete replacement for local search engines

A Federated Search and Retrieval Platform for Lexical Resources in Text+ and CLARIN
eLex 2023: electronic lexicography in the 21st century

10

# CLARIN Federated Content Search (FCS)

A Federated Search and Retrieval Platform for Lexical Resources in Text+ and CLARIN
eLex 2023: electronic lexicography in the 21st century

# Communication Protocol:
# FCS Core 2.0

» Extension of SRU (Search/Retrieval via URL) / searchRetrieve

  » Standardized by Library of Congress LoC / OASIS

  » Data as XML

» RESTful

  » **Explain**: Existing resources

    » Language, annotations, supported data formats, etc.

  » **SearchRetrieve**: search query

» **FCS-QL** as general query language (FCS 2.0)

A Federated Search and Retrieval Platform for Lexical Resources in Text+ and CLARIN
eLex 2023: electronic lexicography in the 21st century

12

# Assumption over data structure

» Full text + optional annotation layers

| Full Text | Die *(The)* | Autos *(cars)* | Sind *(are)* | Schnell *(fast)* |
|---|---|---|---|---|
| Part of Speech (UD17) | DET | NOUN | VERB | ADJ |
| Base form | Das *(The)* | Auto *(car)* | Ist *(is)* | Schnell *(fast)* |
| Phonetic transcription (=SAMPA) | … | … | … | … |
| Orthographic transcription | … | … | … | … |
| Orthographic normalization | … | … | … | … |

A Federated Search and Retrieval Platform for Lexical Resources in Text+ and CLARIN
eLex 2023: electronic lexicography in the 21st century

13

# Query Language FCS-QL

» Similar to CQP (e.g. corpus query workbench)

» Supports multiple annotation layers

A Federated Search and Retrieval Platform for Lexical Resources in Text+ and CLARIN
eLex 2023: electronic lexicography in the 21st century

# Visualization of results (1/2)



[API-Call](#)

# Visualization of results (2/2)

# Existing Resources

» Ecosystem

    » Reference implementations

    » Libraries (Java), Documentation

    » FCS/SRU Validator

    » Endpoint Registry

» Existing resources (actual number fluctuate)

    » 20 institutes in 38 endpoints in 11 countries

    » About 200 „collections" in varying granularity in ~60 languages

A Federated Search and Retrieval Platform for Lexical Resources in Text+ and CLARIN
eLex 2023: electronic lexicography in the 21st century

17

# FCS for Lexical Resources?
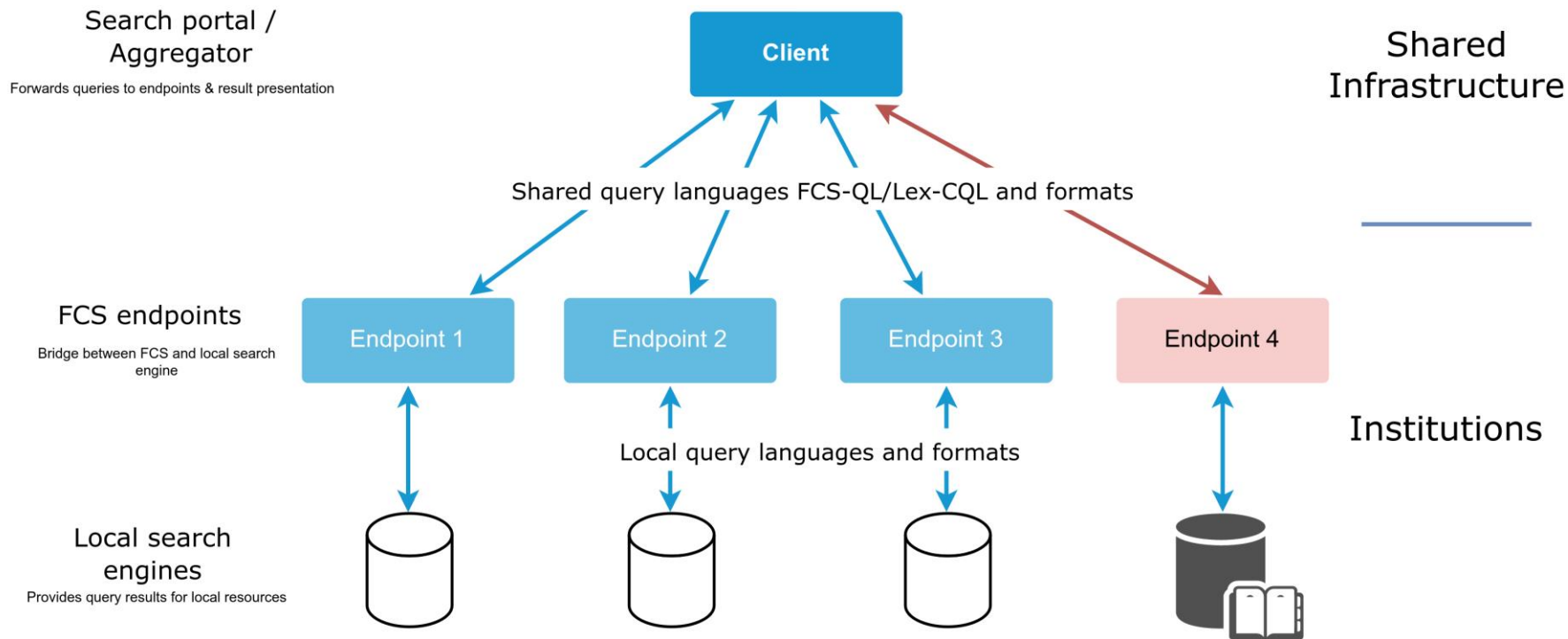
» No – not yet!

» Focus on text „streams" (corpora, transcriptions) & querying of annotation layers

» lexical resources structurally completely different
(word lists, word nets/graphs, key-value based, …)

Specification extension

A Federated Search and Retrieval Platform for Lexical Resources in Text+ and CLARIN
eLex 2023: electronic lexicography in the 21st century

18

# FCS Specification Extension for Lexical Resources Goals

» Query language dedicated to querying lexical entries

  » Subset of **Contextual Query Language** (CQL), agreements on accessible fields of information for a lexeme, complex queries

» Common data formats for unified result presentation

  » Mandatory **LexHITS** data view with **inline annotation** of information types

  » Advanced tabular representation, key-value style

» Compatibility with FCS architecture

  » Reuse of features: access control for restricted resources, automatic registering of endpoints

A Federated Search and Retrieval Platform for Lexical Resources in Text+ and CLARIN
eLex 2023: electronic lexicography in the 21st century

19

# FCS Specification Extension for Lexical Resources



Search portal /
Aggregator

Forwards queries to endpoints & result presentation

Client

Shared
Infrastructure

Shared query languages FCS-QL/Lex-CQL and formats

FCS endpoints

Bridge between FCS and local search engine

Endpoint 1    Endpoint 2    Endpoint 3    Endpoint 4

Institutions

Local query languages and formats

Local search engines

Provides query results for local resources

Text+

A Federated Search and Retrieval Platform for Lexical Resources in Text+ and CLARIN
eLex 2023: electronic lexicography in the 21st century

# FCS Specification Extension for Lexical Resources Draft v0.1

» Query Language – LexCQL

» DataViews – LexHITS + Tabular (draft)

» Draft: https://doi.org/10.5281/zenodo.7849753

A Federated Search and Retrieval Platform for Lexical Resources in Text+ and CLARIN
eLex 2023: electronic lexicography in the 21st century

21

# LexCQL

» Subset of Contextual Query Language (CQL)

» Relation „="

» Operators AND/OR/NOT

» (draft) Fuzzy with „/exact" Modifier

» Fields:

    » lemma,

    » pos (UD17)

    » def

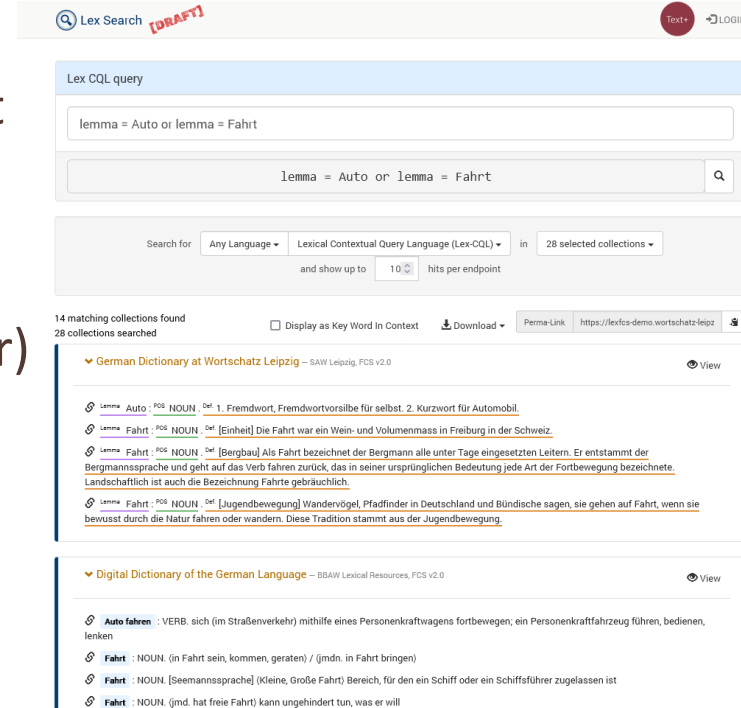    » xr$synonymy, xr$hyponymy, …

    » (draft) senseRef

```
1. cat   # searching on default field, e.g. lemma; specified by endpoint
2. lemma =/exact "läuft"   # exact string match requested
3. def = "an edible" and pos = "NOUN"   # (implicit) partial match in def
4. pos = ADJ and xr$synonymy = "tiny"
5. senseRef = "https://d-nb.info/gnd/118571249"
```

A Federated Search and Retrieval Platform for Lexical Resources in Text+ and CLARIN
eLex 2023: electronic lexicography in the 21st century

22

# LexHITS

» Seamless integration for *mandatory* result format (HITS) BUT optional annotation of „lemma", „pos" and „def"

→ visual hints in frontend (e.g. Aggregator)

```xml
<fcs:DataView type="application/x-textplus-fcs-hits+xml">
  <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>
```

*Apple: NOUN. An apple is an edible fruit produced by an apple tree.*

A Federated Search and Retrieval Platform for Lexical Resources in Text+ and CLARIN
eLex 2023: electronic lexicography in the 21st century

23

# Possible Additional Fields for Complex Tabular Data View

» Structured result presentation, in discussion: „key-value pairs"

» Aim: easy conversion of potential complex formats into general flat structure

» Requires:

  » Recommendation for *required* and *optional* information types

  » *Normative* list of keys and value *formats*

» Attributes:

  » Lemma, Pos, Definition, SenseRef, *nym

  » *Examples, Baseform, Hyphenation, Decomposition*

  » *Cooccurrences, Frequency*

  » *Provenience, Etymology/Word-History*

```xml
<fcs:DataView type="application/x-textplus-fcs-lex+xml">
  <Result>
    <Entry>
      <!-- Lexeme entry -->
      <Name type="lemma">Lemma</Name>
      <Value>Lauf</Value>
    </Entry>
    <Entry>
      <!-- Standard POS tag set -->
      <Name type="pos">POS</Name>
      <!-- Multiple values are possible -->
      <Value>NOUN</Value>
      <Value>VERB</Value>
    </Entry>
    <Entry>
      <!-- Custom POS tag set, as additional "pos" entry type -->
      <Name type="pos">STTS</Name>
      <Value>VVIMP</Value>
      <Value>NN</Value>
    </Entry>
    <!-- ... -->
  </Result>
</fcs:DataView>
```
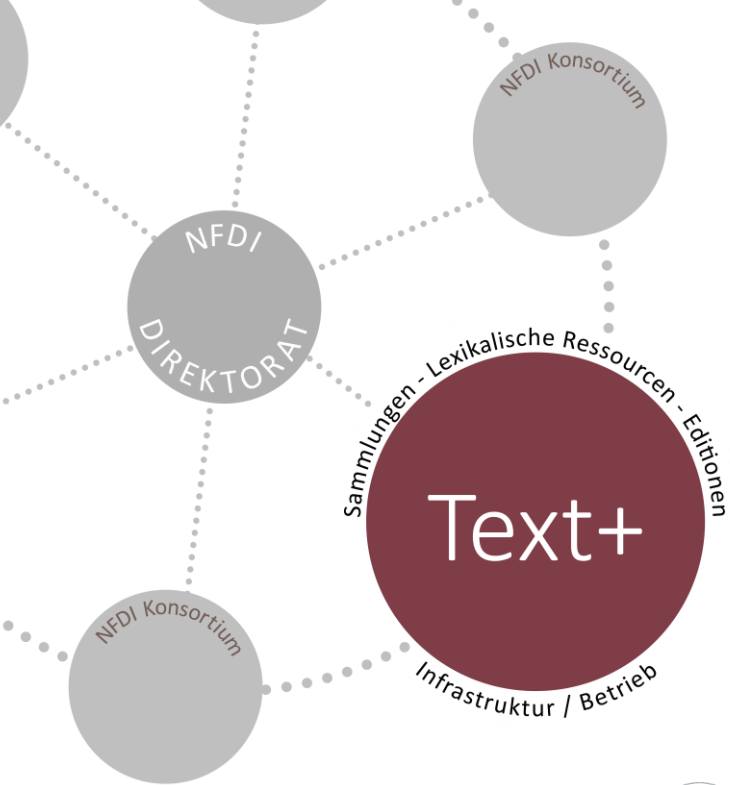
A Federated Search and Retrieval Platform for Lexical Resources in Text+ and CLARIN
eLex 2023: electronic lexicography in the 21st century

24

# Next Steps and Future Work

## Where are we at?

» First draft (v0.1) published
  » LexCQL query language & LexHits data view
  » Demo implementations (Text+ FCS Aggregator + Endpoints)
» Text+: 50 resources from 6 institutes (2 data domains)

## Planned

» Implementation Guide

» Tabular Key-Value Data View

» Endpoint Tester for (Lex)CQL Conformance Levels

A Federated Search and Retrieval Platform for Lexical Resources in Text+ and CLARIN
eLex 2023: electronic lexicography in the 21st century

25

Thank you
for your attention!

Erik Körner, Saxon Academy of Sciences
and Humanities in Leipzig, koerner@saw-leipzig.de

text-plus.org / office@text-plus.org

LexFCS Specification: https://doi.org/10.5281/zenodo.7849753
Text+ LexFCS Aggregator: fcs.text-plus.org
Sources: https://gitlab.gwdg.de/textplus/ag-fcs-lex-fcs-aggregator