Faculty of Environmental Sciences
Chair of Geoinformatics

# Large Language Models for Conversational Geodata Search

AGILE 2025 Tutorial, Dresden, Germany
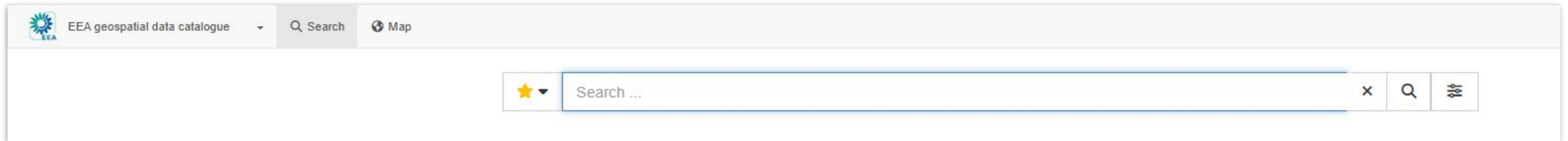
Simeon Wetzel, Auriol Degbelo, Stephan Maes

# Agenda

- **Part 1** (45min)**:**
    - Introduction / Motivation / Scenario

- **Part 2** (75 min)**:**
    - LLM Calls
    - Retrieval Augmented Generation (RAG)
    - Geocoding/ Query Interpretation
    - Conversation
- **Part 3** (75 min)**:**
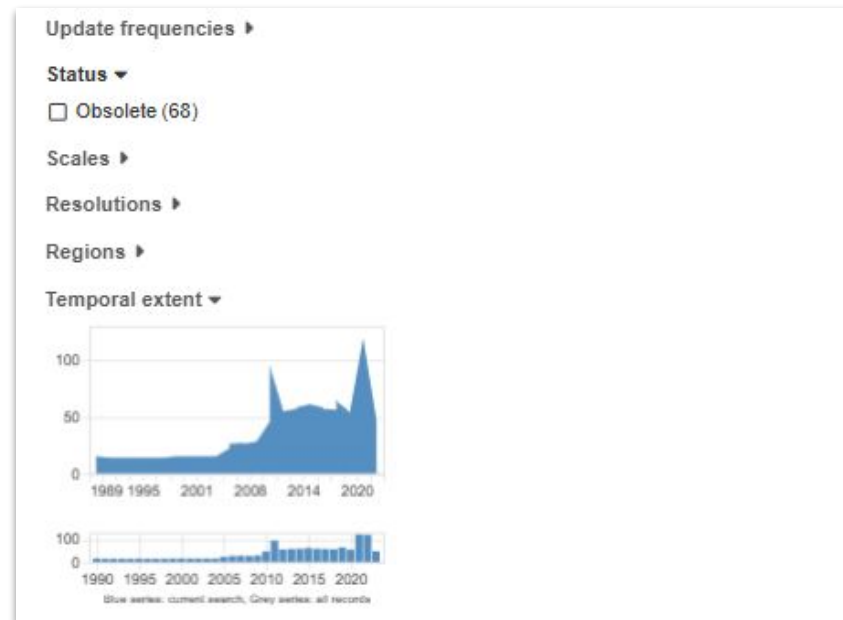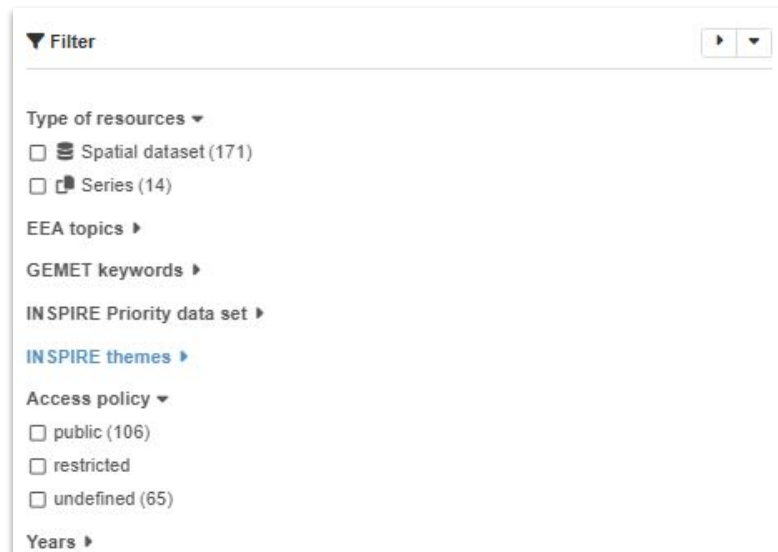    - Agents
    - SmolAgents Framework

# Motivation: LLMs for Conversational Geodata Search

# Traditional Search Approach

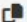- **Metadata Catalogues / Geoportals**  (e.g. EEA SDI Catalogue)
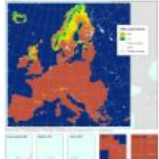    - Full-text interface



    - Search Filter

# Traditional Search Approach

- **Metadata Catalogues / Geoportals**
  - SERP-based result list

# Traditional Search Approach

- **Cycle of single-hop search + refinements**



Full-text search + Filters

Review of the results

Evaluation of „fitness for use"

Query/Filter Refinement

# Challenges with traditional search approach

- Dependency on attributes / metadata quality (completeness / accuracy)

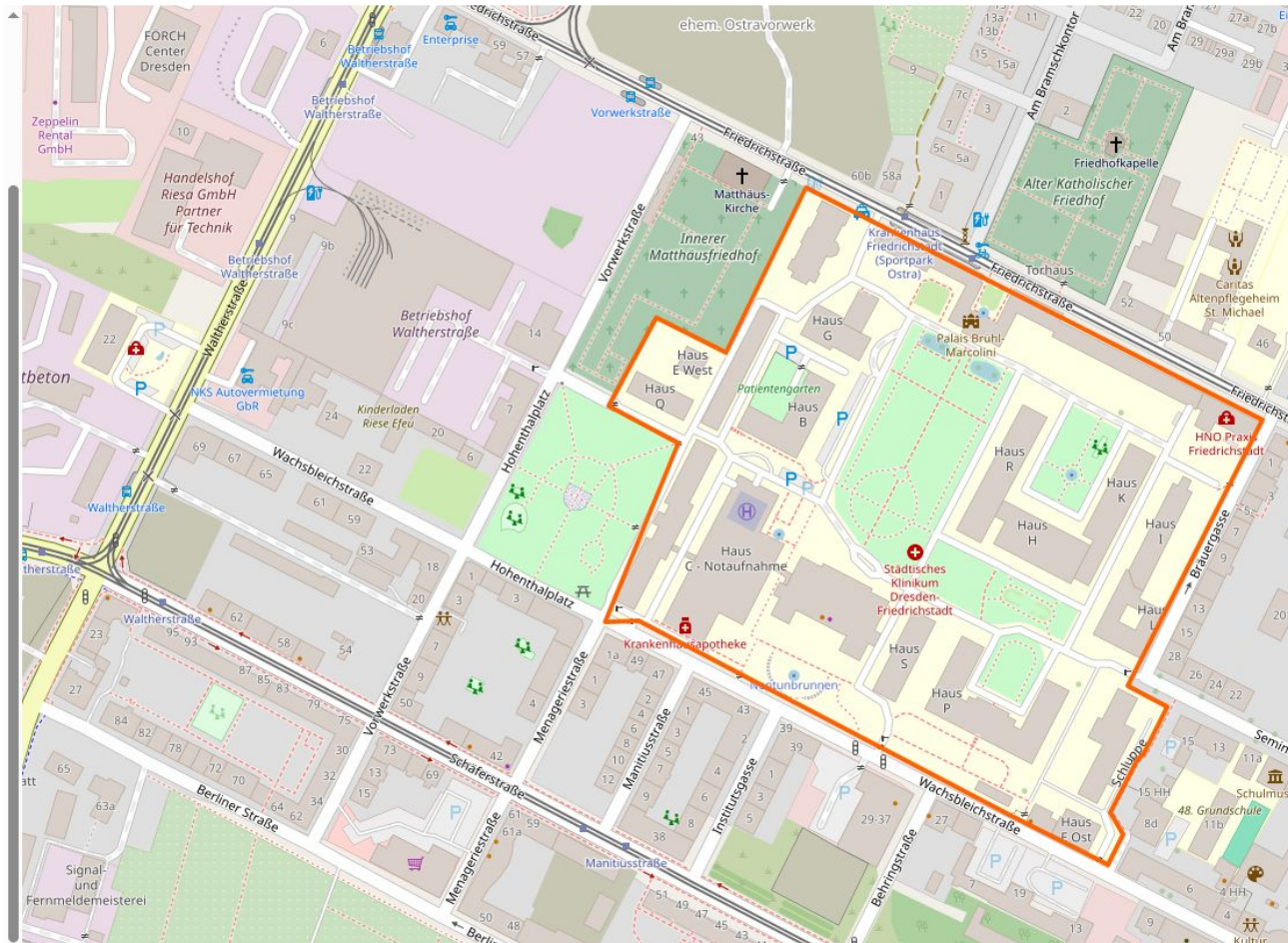# Challenges with traditional search approach

- Dependency on attributes

# Challenges with traditional search approach

- Lexical Search => No semantic search

# Challenges with traditional search approach

- Specific Terminology in ESS data

# Challenges with traditional search approach

- Context-less single-hop queries

- Lexical Search issues

- Specific Terminology in ESS data

- **Complex queries:**

  - 🌍🕐  Spatio-temporal queries:

    - „Historic buildings around…", „"Heavy Precipitation Europe", „Climate Projection 2020-2100"

  - ❓ Ambiguous entities:

    - „Radiation data", „Buildings in Frankfurt Germany"…

  - 🌍🤷  Vague spatial entities:

    - „…East coast…", „North of Ireland"

# Improvements with LLM-based search

- Improved capabilities for …

  … **query interpretation**

  … search **result interpretation**

  … **context-awareness**

  … **semantic search**

# Scenario in this tutorial

TECHNISCHE
UNIVERSITÄT
DRESDEN

DRESDEN
concept

# Scenario: Design of a LLM-driven search architecture for geodata

## Data:

- OpenStreetMap Data (buildings in Dresden, ~50k features used)

## Data Pre-processing

- Representing OSM data as embeddings
- Loading data into a vector store

## LLM-based search:

- Retrieval of the data
- Using data as context

TECHNISCHE
UNIVERSITÄT
DRESDEN

DRESDEN
concept

# Scenario

**Part 1:**

**1. Fetch Data** → **2. Try lexical search**

**3. Prepare data for vector store**

**Two query scenarios**
1. Queries for a **specific buildings** (by name) (e.g. "Deutsches Hygiene Museum")
2. Queries by **building type** (e.g. "museums in Dresden")

# Scenario

**Part 2:**

```
┌──────────────┐        ┌──────────────┐
│  1. Create   │──────▶ │ 2. Issues    │
│  Collection  │        │ with single  │
│              │        │ collection   │
└──────────────┘        └──────────────┘
        │                       │
        ▼                       │
┌──────────────┐        ┌──────────────┐
│ 3. Semantic  │──────▶ │ 4. Query     │
│ Routing      │        │ with multiple│
│              │        │ collections  │
└──────────────┘        └──────────────┘
```
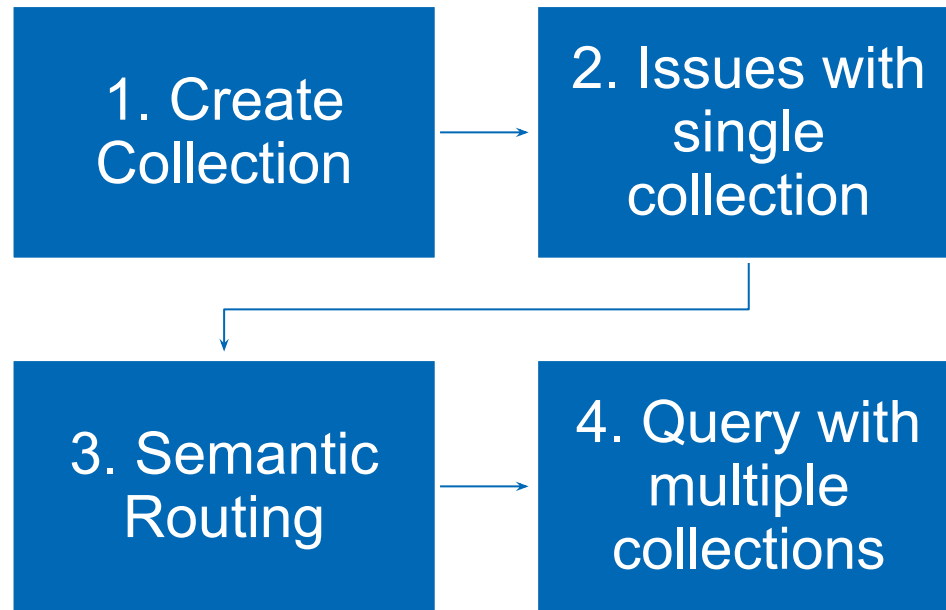
# Notebook:



https://bit.ly/agile25-llm

TECHNISCHE
UNIVERSITÄT
DRESDEN

DRESDEN
concept

# API Keys

https://bit.ly/agile25-key

TECHNISCHE
UNIVERSITÄT
DRESDEN

DRESDEN
concept