

# Developing CLMS Standards for Generative AI Training and Web Crawlers Using Quarto Markdown and Sitemaps

## Task 10.1: Information Provisioning for Generative Chatbots

Ayan Chatterjee, Department of DIGITAL, NILU

2024-09-16

### Index

<b>Abstract</b>	<b>3</b>
<b>1. Introduction</b>	<b>4</b>
1.1. Importance of Copernicus Land Monitoring Service (CLMS) . . . . .	4
1.2. Importance of CLMS Documentation for Web Crawlers: Enhancing Product Discoverability and Findability . . . . .	4
1.3. Web crawling and Information Provisioning for Generative Chatbots . . . . .	6
1.3.1. Motivation . . . . .	7
1.3.2. Importance . . . . .	8
<b>2. Content Standards</b>	<b>8</b>
2.1. Content Structuring . . . . .	9
2.2. HTML Structuring . . . . .	9
2.2.1. Semantic Structuring and Formatting . . . . .	10
2.2.2. Microdata for Enhancing Machine Readability . . . . .	10
2.2.3. Schema Markup for Structured Content . . . . .	10
2.2.4. Headings and Subheadings . . . . .	11
2.2.5. Alt Text and Descriptions . . . . .	11
2.2.6. Meta Tags and Descriptions . . . . .	11
2.2.7. Phrasing and Content Presentation . . . . .	11
2.2.8. Structured Data Repositories . . . . .	12
2.3. PDF Structuring . . . . .	12
2.3.1. Accessible PDF Formats by Tagging . . . . .	13
2.3.2. Structuring and Formatting . . . . .	13

2.3.3. Adding Metadata . . . . .	13
2.3.4. Optimizing Content Presentation . . . . .	14
2.3.5. Setting Up for Knowledge Transfer to Generative AI . . . . .	14
<b>3. Developing CLMS Standards</b>	<b>15</b>
3.1. Tools for CLMS Documentation . . . . .	15
3.2. Quarto Markdown . . . . .	17
3.2.1. Significance . . . . .	17
3.2.2. Configuring Quarto with Jupyter Notebooks . . . . .	18
3.3. Indexing . . . . .	18
3.3.1. Sitemap Generation . . . . .	19
3.3.2. Steps to Implement and Submit the Sitemap . . . . .	20
3.3.3. Enhancing Indexing for Web Crawlers and AI Models . . . . .	20
<b>4. Recommended Standards for Information Formatting</b>	<b>22</b>
4.1. Suggested Standards . . . . .	22
4.2. Guideline for the Process Verification . . . . .	23
<b>5. Conclusion</b>	<b>24</b>
<b>6. References</b>	<b>24</b>

## Abstract

Generative chatbots rely on large amounts of structured data to provide accurate, timely responses to user queries. By developing **Copernicus Land Monitoring Service (CLMS)** standards for information formatting and delivery using **Quarto Markdown** and **sitemaps**, we can ensure that the vast amounts of environmental data in CLMS are accessible to web crawlers and AI models. Using standardized structured content improves discoverability and discoverability of CLMS products and makes it easier for users to access relevant datasets through traditional search engines and generative chatbots.

In addition, by providing clear guidelines for content formatting, cross-referencing, and sitemap management, this approach ensures that the CLMS data repository remains up-to-date and well-organized. This in turn supports the training of AI models to help users find exactly the CLMS products they need, whether through direct query or generative chatbot interaction.

# 1. Introduction

## 1.1. Importance of Copernicus Land Monitoring Service (CLMS)

The Copernicus Land Monitoring Service (CLMS) is a critical component of the Copernicus Programme, which is the European Union’s Earth observation initiative [1]. The service is responsible for providing timely and accurate land cover and land use data, along with a wide range of environmental variables related to land ecosystems. This data is essential for understanding and managing Europe’s environmental resources, supporting sustainable development, climate monitoring, and informed policy-making. The key areas where CLMS is vital include:

- **Environmental Monitoring:** CLMS provides data on land cover, vegetation, soil, and water bodies, which are crucial for monitoring environmental changes such as deforestation, urban sprawl, and the health of ecosystems. This data supports conservation efforts and helps in tracking biodiversity and land degradation.
- **Sustainable Land Management:** With the growing need for sustainable practices, CLMS delivers data that helps governments and organizations plan and manage land resources more effectively. It supports agriculture, forestry, water management, and urban planning, helping to mitigate the effects of climate change.
- **Climate Change Monitoring:** CLMS plays a significant role in assessing the impact of climate change on European landscapes. It helps track changes in land use, vegetation, and land surface temperatures, which are important indicators of climate change impacts.
- **Disaster Management:** CLMS data is used for emergency response and disaster management, especially in cases of floods, fires, and other natural disasters. The accurate and near-real-time data allows authorities to take preventive actions and make quick decisions during emergencies.
- **Policy Support and Decision-Making:** The service supports EU environmental policies, including the Green Deal, Common Agricultural Policy (CAP), and the EU Biodiversity Strategy. The data provided by CLMS informs decision-makers at the European, national, and local levels, ensuring that policies are grounded in the latest environmental data.

## 1.2. Importance of CLMS Documentation for Web Crawlers: Enhancing Product Discoverability and Findability

The discoverability and findability of CLMS products on the web are crucial for ensuring that this valuable environmental data is accessible to a wide range of users, including researchers, policymakers, and environmental organizations. Making CLMS documentation available on

the web for crawlers facilitates product discoverability by enabling search engines and AI-powered systems (like generative chatbots) to index, retrieve, and present relevant data to users. Here's why ensuring that CLMS documents are available to web crawlers is essential:

- **Increased Accessibility for Diverse Users:** CLMS products cater to a broad audience, including government agencies, NGOs, scientists, and the public. Properly formatted and exposed documentation allows these users to easily find and access data via search engines. Web crawlers can efficiently index CLMS products, simplifying the search for specific datasets without navigating complex databases.
- **Enhanced Search Engine Optimization (SEO):** CLMS products cater to a broad audience, including government agencies, NGOs, scientists, and the public. Properly formatted and exposed documentation allows these users to easily find and access data via search engines. Web crawlers can efficiently index CLMS products, simplifying the search for specific datasets without navigating complex databases.
- **Improved Product Findability Through AI and Chatbots:** AI-powered search tools and chatbots use indexed information to generate responses. By ensuring that CLMS documentation is structured for crawling, CLMS products become accessible to third-party chatbots, expanding their reach through natural language queries and conversational interfaces.
- **Faster and More Accurate Data Retrieval:** Well-formatted CLMS documents enable faster and more accurate data retrieval, essential for time-sensitive applications like disaster management. Proper crawling ensures that search engines and AI systems provide up-to-date CLMS products, crucial for timely decision-making.
- **Standardization and Interoperability:** Adopting CLMS standards and formats like Quarto Markdown ensures consistency, making documents easier to index and retrieve. Standardization promotes interoperability, allowing CLMS data to be used across various platforms, including AI systems and environmental tools.
- **Global Reach and Broader Impact:** Making CLMS documents available to web crawlers increases their global reach. Optimized data allows users worldwide to access key environmental information, contributing to global initiatives, research, and policymaking beyond the EU.
- **Supporting Third-Party Integration:** Third-party platforms rely on web crawlers and AI tools to access CLMS data. By exposing CLMS products to crawlers, the data can be integrated into various tools and services, enhancing discoverability and promoting broader use in AI-driven analytics and public services.

By making CLMS documents available to web crawlers using standardized formats such as HTML, PDF, and DOCX (which adhere to semantic structure, web standards, and use meta-data), CLMS can ensure that its products are easily indexed, retrieved, and integrated into a variety of search engines, artificial intelligence systems, and chatbots. This strategy not only

increases the visibility of CLMS products, but also improves accessibility to a global audience, ensuring that researchers, policymakers, and the public can effectively find and use CLMS data. At a time when timely, accurate environmental data is becoming increasingly important, optimizing CLMS products for web crawlers is a necessary step to ensure that everyone has access to these valuable resources.

### 1.3. Web crawling and Information Provisioning for Generative Chatbots

Web crawling is the process used by search engines to explore and index the web pages of websites. The crawler downloads pages, reads the content, and adds it to the search engine's index. Crawlers are designed to navigate from one page to another by following hyperlinks, allowing them to efficiently cover a website's entire structure. Search engines rely on crawlers to keep their results up-to-date by regularly visiting websites and checking for new or modified content. **Googlebot**, **Bingbot**, and **Yahoo Slurp** are some example of popular web crawlers. Key terms involved in web crawling are:

- **Search engine:** A system that allows users to search for content on the web.
- **Indexing:** The process of storing web content so it can be retrieved later.
- **Web pages:** Documents that make up the web, interconnected by hyperlinks.
- **Hyperlinks:** Links that connect different web pages, forming a navigable web.

Web crawling has become essential for search engines and AI applications. The integration of these technologies has been explored extensively [2], [3], [4], [5]. The growth of digital content has placed significant demands on the efficiency and accuracy of web crawlers and artificial intelligence (AI) models [6], [7]. In response, Content Lifecycle Management Standards (CLMS) are essential for establishing uniformity in the way data is formatted, structured, and exposed for automated tools like crawlers and AI training datasets. CLMS helps ensure that content is easy to access, interpret, and process, leading to more accurate information retrieval and AI model training. This document outlines the development of CLMS standards for exposing information to web crawlers and optimizing the formatting for AI data ingestion. Figure 1 focuses on the working of a web crawler [8].

In recent years, generative chatbots have made great progress and become powerful tools that allow users to access detailed information and conduct complex queries. In particular, chatbots can help users explore certain aspects of CLMS products, such as allocation rules or the purpose of a particular product. These tools are not only critical for product discoverability, but also improve user understanding of CLMS products. To ensure that chatbots effectively help users find and understand CLMS products, it is important that the underlying information is formatted and presented in a way that is easy to find and use. This requires well-structured documentation and a system that allows web crawlers and AI models to effectively access and process CLMS data.

Web crawlers and AI models are critical to the discoverability of online information. Web crawlers that index websites rely on well-structured content to perform their tasks effectively.

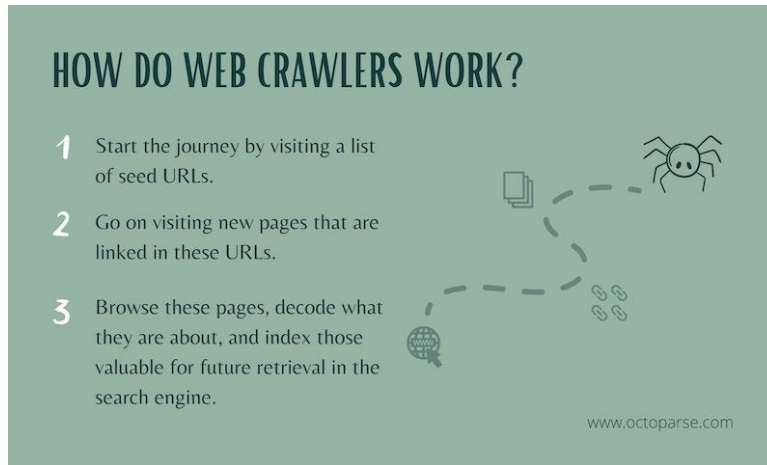


Figure 1: Diagram illustrating web crawling [8].

Similarly, generative AI models, including chatbots, require high-quality structured data to produce accurate and meaningful results. CLMS provides important environmental data, but in order for this data to be useful to AI models and easy for users to find, it must be properly formatted and made available.

### 1.3.1. Motivation

The relationship between AI and web crawlers has led to new frontiers in both industries. The primary motivation for creating CLMS standards lies in the need for:

- **Improved Crawling Efficiency:** Properly formatted content with metadata helps crawlers index relevant information faster and more accurately.
- **Better AI Model Training:** Consistent content structure ensures that AI models are trained on high-quality, organized data.
- **Data Accessibility:** Standardizing the structure of content ensures that information is universally accessible across platforms.

The following key aspects are critical for ensuring that data is structured and accessible for web crawlers and AI systems:

- **Uniform metadata:** Consistent metadata usage across all content is essential. Metadata includes details like title, author, keywords, and publication date. Uniform metadata ensures that web crawlers and AI systems can easily index and categorize content, improving searchability and discoverability.

- **Clearly defined content sections:** Content should be organized into distinct sections, such as titles, headings, and subheadings. This structured format helps both users and machines navigate through the content efficiently, making key information easy to locate and retrieve.
- **Embedded structured data formats:** Incorporating structured data formats such as **JSON-LD**, **RDF**, or **XML** provides a precise way of representing information. These formats help web crawlers and AI systems understand relationships and attributes within the content, facilitating accurate extraction, interpretation, and use of the data across various platforms.

### 1.3.2. Importance

- **Enhanced Web Crawling:** Properly structured CLMS content will improve web crawlers' ability to index and retrieve information.
- **Improved AI Training:** Structured data ensures higher-quality datasets, which result in better-trained AI models, particularly for generative chatbots.
- **Better User Experience:** By improving product discoverability and findability, users will have an easier time accessing and understanding CLMS products.

#### Tip

Given the growing complexity of CLMS products and the increasing reliance on generative AI tools, it is critical to implement standards that improve the discoverability and usability of CLMS data.

#### Note

By standardizing the format and delivery of CLMS information, our goal is to ensure that generative AI applications, such as web crawlers and chatbots, can effectively access and use this data.

## 2. Content Standards

Developing content standards requires collaboration between content creators, data engineers, and AI researchers. The process typically follows these stages for different document types in use:



## 2.1. Content Structuring

Content structuring involves organizing data into recognizable, standard components, such as:

- **Title:** Main identifier of the content.
- **Metadata:** Information about the content, including authors, dates, keywords, and relevant classification.
- **Headings and Subheadings:** Structured sections that break down the content into digestible parts.

The example of Metadata formatting has been given below:

```
---
title: "Developing CLMS Standards for Generative AI Training and Web Crawlers"
subtitle: "Task 10.1: Information Provisioning for Generative Chatbots"
author: "Ayan Chatterjee, Department of DIGITAL, NILU, ayan@nilu.no."
date: "2024-09-10"
sitemap: true           #Enables sitemap generation for web crawlers
toc: true               # Enable the Table of Contents
toc-title: "Index"      # Customize the title of the table of contents
toc-depth: 3           # Include headings up to level 3 (###)
keywords: ["CLMS standards", "web crawlers", "AI training", "information formatting"]
bibliography: references.bib # Link to the bibliography file
csl: ieee.csl           # Link to the CSL file for IEEE style
format:
  html: default
  pdf: default
  docx: default
---
```

## 2.2. HTML Structuring

The following structured approach in HTML allows web crawlers to effectively index and retrieve content while facilitating AI training for generative models, ensuring that information is both accessible and usable:

### 2.2.1. Semantic Structuring and Formatting

It is used to enhance both **machine readability** and **user comprehension**, we must follow structured and semantic formatting principles. This includes using HTML5 elements, schema markup, and providing clear metadata. Using HTML5 semantic elements like `<article>`, `<section>`, `<header>`, and `<footer>` helps structure the document meaningfully. For example:

```
<article>
  <header>
    <h1>Understanding Web Crawlers</h1>
    <meta name="description" content="How web crawlers work and index ..!" />
  </header>
  <section>
    <h2>How Crawlers Index Content</h2>
    <p>Web crawlers use semantic structure to efficiently index web pages.</p>
  </section>
  <footer>
    <p>Author: Ayan Chatterjee</p>
  </footer>
</article>
```

### 2.2.2. Microdata for Enhancing Machine Readability

Microdata attributes such as `itemscope`, `itemtype`, and `itemprop` provide semantic clarity for machines, enabling more efficient crawling and interpretation.

```
<article itemscope itemtype="https://schema.org/Article">
  <header>
    <h1 itemprop="headline">Web Crawling Explained</h1>
    <meta itemprop="description" content="How web crawlers index ..?" />
  </header>
</article>
```

### 2.2.3. Schema Markup for Structured Content

Use Schema Markup (like `ResearchArticle`, `Dataset`, or `CreativeWork`) to define the content type and enhance machine readability. This helps both web crawlers and AI to categorize content accurately.

```
<article itemscope itemtype="https://schema.org/ResearchArticle">
  <header>
    <h1 itemprop="headline">AI Training for Web Crawlers</h1>
    <meta itemprop="description" content=" AI training techniques for .." />
  </header>
</article>
```

#### 2.2.4. Headings and Subheadings

Provide clearly defined headings and subheadings to organize content for easier navigation and indexing by crawlers.

```
---
# How AI Models are Trained
## Data Collection
## Model Training
## Evaluation
---
```

#### 2.2.5. Alt Text and Descriptions

For images and diagrams, always provide alt text and descriptions to improve accessibility.

```
![A diagram illustrating how web crawlers work]
(images/web_crawlers.png){alt="A diagram of web crawler processes" width=50%}
```

#### 2.2.6. Meta Tags and Descriptions

Add meta tags and descriptions to help web crawlers index the content more accurately

```
<meta name="description" content="How web crawlers work effectively!" />
```

#### 2.2.7. Phrasing and Content Presentation

Ensure that important keywords are present in titles, headings, and throughout the content without overusing them (avoid keyword stuffing).

```
# Introduction to Web Crawlers and AI Training
Web crawlers, also known as spiders, are used by search engines to index web ...
```

Write in a clear and concise manner. Avoid jargon unless necessary, and ensure that key concepts are easy to understand.

```
Web crawlers automatically scan websites to collect and index content.
They follow links, downloading web pages and saving them for future queries.
```

Use hyperlinks and cross-references to guide both users and web crawlers to related content.

```
For more details, see the [Introduction to AI Training](#data-collection).
```

Provide a brief abstract or summary at the beginning of each article or section for better clarity and indexing.

```
**Summary:** This article provides an overview of indexing content,
and their integration with AI.
```

### 2.2.8. Structured Data Repositories

It is used to enable knowledge transfer to generative AI, use standardized formats like JSON-LD, RDF, or XML to define metadata and structure.

```
{
  "@context": "https://schema.org",
  "@type": "Dataset",
  "name": "AI Training Dataset",
  "description": "A dataset designed to improve search engine crawlers."
}
```

```
<dataset xmlns="http://www.w3.org/2001/XMLSchema-instance" type="AI Training Dataset">
  <name>AI Training Dataset</name>
  <description>A dataset designed for training AI models.</description>
</dataset>
```

### 2.3. PDF Structuring

The following structured approach in PDF will improve documents for indexing by web crawlers, integration with AI systems, and overall improved accessibility for users:

### 2.3.1. Accessible PDF Formats by Tagging

Ensure that the PDF is tagged properly so that screen readers and AI tools can interpret the document structure. For instance, headings, paragraphs, and lists should be tagged semantically.

```
# Heading 1 (tagged as <h1>)
- List item 1 (tagged as <ul><li>)
```

### 2.3.2. Structuring and Formatting

The document structure should be accessible, with a clear hierarchy and a clickable table of contents (TOC). Accessible tagging, hierarchical organization, and text over image improve the usability for both humans and machines.

Organize content into a well-defined hierarchy using headings (#, ##, ###). This improves both user navigation and machine parsing for AI and web crawlers.

```
## Section 1: Introduction
### Subsection 1.1: Overview

toc: true
toc-depth: 2
```

### 2.3.3. Adding Metadata

Embedding metadata such as document properties (e.g., Title, Author, Subject, and Keywords), XMP metadata, Schema.org metadata, and descriptive metadata helps search engines and AI systems index, categorize, and retrieve information efficiently.

```
title: "PDF Structuring and Formatting"
author: "Ayan Chatterjee"
subject: "Document Accessibility and Metadata"
keywords: ["PDF accessibility", "metadata", "AI integration"]
```

XMP metadata is stored as XML in the PDF file, allowing for rich data descriptions. Schema.org metadata in JSON-LD provide structured information that AI and web crawlers can easily understand.

```
{
  "@context": "https://schema.org",
  "@type": "CreativeWork",
  "name": "PDF Structuring and Formatting",
  "author": {
    "@type": "Person",
    "name": "Jane Doe"
  },
  "keywords": ["PDF accessibility", "metadata", "AI integration"]
}
```

### 2.3.4. Optimizing Content Presentation

Ensuring the proper placement of keywords, providing alt text for images, and correctly labeling figures and tables contribute to the searchability and accessibility of the content. This is crucial for effective interaction with web crawlers and AI models.

```
Keywords: PDF accessibility, web crawlers, generative AI
! [A flowchart showing the PDF processing workflow] (path/to/image.png){alt="PDF workflow"}
! [Figure 1: A table of contents structure] (path/to/image.png){#fig-toc}
```

### 2.3.5. Setting Up for Knowledge Transfer to Generative AI

Using machine-readable fonts (e.g., Arial, Times New Roman), a clean and simple layout, and adding comments or annotations helps prepare the document for use in generative AI systems. AI models benefit from well-structured and easy-to-parse content, which improves their ability to understand and generate meaningful responses based on the content.

```
## Section 1: Overview
This section introduces the importance of accessible PDFs for AI processing...

<!-- This annotation explains the role of hierarchical metadata for AI -->
```

#### ! Important

By such structured practices, we can ensure that the content is both human-readable and machine-readable, facilitating easy discovery by web crawlers and seamless integration with AI training systems.

### 3. Developing CLMS Standards

In the context of **Developing CLMS Standards**, it is essential to utilize advanced tools that support both the creation of well-structured documents and the easy discoverability of content for web crawlers and AI systems. Several tools are available for content formatting, documentation, and publication. Among these, **Quarto** stands out due to its versatility, allowing users to create, format, and publish documents in multiple formats (HTML, PDF, Word) with integrated support for code execution and structured content.

This section compares several of these tools, explaining why **Quarto** is particularly suitable for creating CLMS-compliant documentation. We'll also cover how to configure Quarto with **Jupyter Notebooks** and the importance of using **Quarto Markdown** for CLMS content. A Quarto Markdown file provides a structured approach to documenting the development of CLMS standards, ensuring content is easily accessible by both web crawlers and AI systems.

#### 3.1. Tools for CLMS Documentation

- **Quarto:** Quarto is a highly versatile tool for creating and publishing documents, including PDFs, with rich formatting, code integration, and support for multiple formats (HTML, PDF, Word). Quarto's cross-platform capabilities make it ideal for creating structured and searchable documents for CLMS, supporting web crawlers and AI applications.
- **R Markdown:** A popular tool in the R community that allows users to combine narrative text with R code, producing output in HTML, PDF, and Word formats. Though powerful for statistical analysis, it is more limited in non-R-based workflows compared to Quarto.
- **Jupyter Notebooks:** An interactive tool supporting over 40 programming languages, commonly used for data science and computing. Notebooks can be exported to multiple formats (HTML, PDF, slides), but lack Quarto's advanced content formatting features.
- **Pandoc:** A universal document converter that enables conversion between various markup formats, including Markdown, LaTeX, and HTML. While powerful for conversions, Pandoc lacks the code integration and dynamic formatting of Quarto.
- **LaTeX:** A document preparation system for producing scientific and technical documents. While highly customizable, it requires significant expertise and lacks the ease of Markdown tools like Quarto.
- **Hugo:** A static site generator used for creating websites and blogs from Markdown files. While efficient for websites, it doesn't provide the same level of document control and integration as Quarto.

- **Sphinx:** A documentation generator mainly used for Python projects. It supports conversion to formats like HTML and PDF but lacks the cross-language support and document versatility of Quarto.
- **Bookdown:** An extension of R Markdown, designed for writing books and long documents. It supports multiple output formats but is mostly R-focused, while Quarto supports multiple languages.
- **GitBook:** A tool for creating documentation and books using Markdown. It allows collaboration but lacks the dynamic formatting and multi-language support found in Quarto.
- **Pelican:** A static site generator that uses Markdown or reStructuredText. Best suited for blogs, it doesn't provide the integrated support for complex documents required by CLMS standards.
- **Typora:** A WYSIWYG Markdown editor that offers easy editing but lacks the advanced document control and integration capabilities that Quarto provides.

The comparison of tools for CLMS documentation as shown in below Table 1. As shown in Table 1, Quarto outperforms other tools in terms of supported output formats and reproducibility.

Table 1: Comparative analysis of Quarto versus other formatting tools.

Tool	Cross-Language Support	Output Formats	Code Integration	Static Site Generation	Ideal Use Case
<b>Quarto</b>	Yes	HTML, PDF, Word	Yes	Yes	Reports, blogs, CLMS docs
R Markdown	R only	HTML, PDF, Word	Yes (R)	No	Statistical reports
Jupyter Notebooks	40+ languages	HTML, PDF	Yes	No	Data Science
LaTeX	Limited	PDF, HTML	No	No	Scientific papers
Hugo	No	HTML	No	Yes	Blogs, websites
Sphinx	Python	HTML, PDF	No	Yes	Python documentation



## 3.2. Quarto Markdown

Markdown is a lightweight, easy-to-read syntax used for formatting plain text documents [9], [10], [11]. In Quarto, Markdown is extended to support additional features beyond standard Markdown, allowing users to write text, integrate code, and generate richly formatted documents in various formats such as HTML, PDF, and Word [9], [10], [11]. Quarto Markdown combines the simplicity of regular Markdown with powerful features for document rendering, making it ideal for data analysis, technical writing, academic papers, and reports [9], [10], [11].

Quarto Markdown uses the standard Markdown syntax for headings, lists, emphasis, and links, while also supporting enhanced features like cross-referencing, citations, figures, tables, mathematical equations, and more [9], [10], [11]. Quarto also allows for code execution in multiple programming languages (such as Python, R, and Julia) embedded within the Markdown file, enabling dynamic document creation where the outputs are generated directly from the code [9], [10], [11], [12].

Key features of **Markdown** in **Quarto** are:

- **Standard Markdown:** Supports headings, lists, links, images, bold, italics, etc.
- **YAML Header:** Allows users to specify metadata like title, author, date, and output formats (HTML, PDF, Word) at the start of the document.
- **Cross-references:** Provides automatic numbering and referencing for figures, tables, sections, etc.
- **Code Execution:** Integrates code cells for multiple programming languages, making it possible to run code and include its outputs directly in the document.
- **Mathematics and Equations:** Supports LaTeX-style equations for technical writing.
- **Citations:** Allows for referencing research papers and articles using BibTeX or CSL styles.
- **Multi-output Format:** Enables seamless conversion to multiple formats like HTML, PDF, Word, presentations, and slides.

### 3.2.1. Significance

Markdown in Quarto can be significant due to its **simplicity and flexibility** for CLMS documentation. With an **easy-to-use syntax**, it allows users to format text without requiring complex tools, making it accessible to both non-technical users and programmers. This flexibility enables the creation of a wide variety of documents, ranging from blog posts to scientific reports. Quarto extends standard Markdown by supporting **rich formatting options** essential for technical and academic writing, including built-in support for tables, figures, equations, footnotes, and cross-referencing. The **integration of code and text** is another powerful feature, allowing Quarto Markdown to embed code execution within documents. This is critical for reproducible research, enabling the inclusion of tables, charts, and figures generated directly

from code, making it highly suitable for data science and technical reporting. Additionally, Quarto Markdown supports **multi-format output**, allowing users to create content once and export it to multiple formats like HTML, PDF, and Word, streamlining document preparation for different audiences. When used for online content, its structured format **improves SEO (Search Engine Optimization)**, making it easier for search engines to index and enhance discoverability. The ease of **managing references, citations, and cross-references** further strengthens its utility in academic and research documentation. Since Markdown files are plain text, Quarto seamlessly integrates with **version control** tools like Git, enabling easy **collaboration** among multiple contributors, especially in open-source and research communities. Finally, Quarto Markdown’s versatility in document creation extends across blogs, technical documentation, reports, scientific papers, and books, making it an ideal tool for content creators across various disciplines.

### 3.2.2. Configuring Quarto with Jupyter Notebooks

To integrate **Quarto** with **Jupyter Notebooks**:

- **Install Quarto:** Download and install Quarto from [Quarto.org](https://quarto.org).
- **Install Jupyter:** Ensure you have Jupyter installed. If not, install it using `pip`: “`bash pip install notebook`”
- **Rendering:** You can directly write your content in Jupyter Notebooks and then render the notebook using Quarto to multiple formats: “`bash quarto render your-notebook.ipynb -to html quarto render your-notebook.ipynb -to pdf quarto render your-notebook.ipynb -to docx`”
- **YAML Header in Jupyter:**

```
---  
title: "CLMS Data Analysis"  
author: "Ayan Chatterjee"  
format:  
  html: default  
  pdf: default  
  docx: default  
---
```

### 3.3. Indexing

Proper indexing is essential for increasing the discoverability and accessibility of CLMS products [13], [14]. By formatting documents using Quarto Markdown and generating a `sitemap.xml`, we can ensure that search engines and AI systems efficiently crawl and retrieve

CLMS content [13], [14]. To improve document indexing for enhanced discoverability and accessibility we can adopt the following approaches:

- Organize content using **structured headers** and **metadata** in Quarto Markdown.
- Use proper keywords and descriptions in the document metadata.
- Cross-reference related documents to create interconnected content that helps crawlers navigate.

```
---
title: "Land Use Mapping with CLMS Data"
author: "Ayan Chatterjee"
date: "2024-08-01"
keywords: ["land use", "CLMS", "mapping", "environment"]
description: "A detailed report on how CLMS data."
---
```

### 3.3.1. Sitemap Generation

A sitemap.xml helps web crawlers discover all the content on the website [13], [14]. By providing a clear roadmap, crawlers can index each document, ensuring that all CLMS resources are available for search and AI training. By using **Quarto Markdown** and generating a **sitemap.xml**, CLMS documents can be structured in a way that improves their **indexing**, making them more **discoverable** by search engines and AI systems. This approach ensures efficient crawling, improves search engine ranking, and enhances the accessibility of CLMS products for users and AI models alike.

- **Search Engine Discoverability:** Users and AI systems can easily find the indexed CLMS documents.
- **Efficient Crawling:** The sitemap provides a roadmap, allowing for faster and more accurate indexing.
- **Increased Accessibility:** Properly indexed documents are easier for users and AI to retrieve and utilize, improving the overall product visibility.

```
<urlset xmlns="http://www.sitemaps.org/schemas/sitemap/0.9">
  <url>
    <loc>http://example.com/clms/land-use-mapping</loc>
    <lastmod>2024-08-01</lastmod>
    <changefreq>monthly</changefreq>
  </url>
  <url>
    <loc>http://example.com/clms/land-cover-change</loc>
    <lastmod>2024-07-15</lastmod>
```

```
<changefreq>monthly</changefreq>
</url>
</urlset>
```

### 3.3.2. Steps to Implement and Submit the Sitemap

- **Generate the Sitemap:** Use a sitemap generator tool (e.g., XML-Sitemaps or Screaming Frog) to create a sitemap, or have it generated automatically by a CMS like WordPress or a static site generator like Hugo.
- **Upload the Sitemap:** Once generated, place the sitemap.xml file in the root directory of your website, e.g., <https://www.example.com/sitemap.xml>.
- **Submit to Search Engines:** Submit your sitemap to search engines via tools like Google Search Console and Bing Webmaster Tools. This helps search engines index your site properly.

### 3.3.3. Enhancing Indexing for Web Crawlers and AI Models

To ensure that CLMS documents are findable and accessible to web crawlers and AI models, it's important to implement proper steps for generating and submitting a sitemap and using structured data (such as metadata and JSON-LD) to enhance indexing.

- **Descriptive Filenames:** Use filenames that clearly describe the content of the document. For instance, instead of doc1.md, use clms-land-monitoring-data.md.
- **Metadata:** Add descriptive metadata in your Quarto Markdown files (e.g., title, author, keywords). This helps search engines and AI models understand the content better.
- **Text Content:** Ensure that text content is descriptive and structured using headings and subheadings to guide crawlers.
- **HTML Metadata and JSON-LD Structured Data:** Use HTML metadata and JSON-LD structured data within the Quarto document to improve how your content is indexed by search engines and used by AI training systems.

The following Quarto Markdown YAML header example demonstrates how to enhance document visibility for web crawling and AI training by including metadata and structured data. This can be part of your CLMS documentation to ensure that it is well-indexed and easy to discover.

```

---
title: "CLMS Land Monitoring Data"
author: "Ayan Chatterjee"
date: "2024-09-15"
keywords: ["CLMS", "web crawling", "AI training", "environmental data"]
description: "Comprehensive overview of CLMS land monitoring datasets, ....."
sitemap: true # Flag to include this document in the sitemap

# HTML metadata for SEO and discoverability
meta:
  - name: "description"
    content: "CLMS land monitoring datasets for environmental and climate ..."
  - name: "keywords"
    content: "CLMS, land monitoring, environmental data, AI, web crawling"

# JSON-LD structured data to help search engines and AI understand the content
json-ld:
  - "@context": "https://schema.org"
    "@type": "Dataset"
    "name": "CLMS Land Monitoring Data"
    "description": "Detailed data on land monitoring and ...."
    "url": "https://www.example.com/clms-land-monitoring-data"
    "keywords": "land monitoring, environmental data, AI training.."
    "datePublished": "2024-09-15"
    "creator":
      "@type": "Organization"
      "name": "Copernicus Land Monitoring Service"
    "publisher":
      "@type": "Organization"
      "name": "European Environment Agency"

```

### ! Important

**Quarto** stands out as the most versatile tool for creating CLMS-compliant documents, with cross-language support, integration of code, multiple output formats, and the ability to generate static websites.

### ! Important

To ensure that CLMS documents are findable and accessible to web crawlers and AI models, it's important to implement proper steps for generating and submitting a sitemap and using structured data (such as metadata and JSON-LD) to enhance indexing.

## 4. Recommended Standards for Information Formatting

### 4.1. Suggested Standards

One of the main challenges in this task is improving the findability and discoverability of CLMS products. With the extensive range of data and services offered by CLMS, users often struggle to locate specific datasets or resources. Chatbots serve as a potential solution by guiding users to the appropriate resources. For chatbots to effectively perform this role, the data must be properly structured, categorized, and indexed. To support this:

- Documentation must be **accessible to third-party chatbots**. While CLMS chatbots will be the primary interaction point, external platforms should also access and retrieve relevant data. Exposing CLMS data in a structured and standardized format ensures interoperability across various chatbot systems, enhancing discoverability.
- Recommendations will be provided on how CLMS should **format and expose information**. These guidelines will focus on best practices for metadata structuring, content organization, and linkable resources to optimize data formatting.

The recommended standards for CLMS will include the use of Quarto Markdown, sitemaps, and structured metadata for web crawlers and AI systems.

- **Using Quarto Markdown for Data Structuring:** Quarto Markdown allows for the clear organization of data, with structured sections such as headings, subheadings, and metadata fields. This makes it easier for web crawlers and AI systems to navigate the content and retrieve relevant information. Additionally, by using cross-referencing within Quarto Markdown documents, CLMS products and resources can be interconnected, providing users with a more seamless exploration experience.
- **Implementing Sitemaps for Efficient Crawling:** Sitemaps provide a roadmap for web crawlers, ensuring that all relevant pages and data sources are indexed. By creating comprehensive sitemaps that expose the entirety of the CLMS data repository, the task ensures that web crawlers and AI systems can efficiently discover and retrieve content. This is essential for making CLMS data easily accessible to third-party chatbots and AI platforms.

## 4.2. Guideline for the Process Verification

We can compare the results of the search queries for both the unformatted and formatted documents. Typically, formatted documents with clear structure and metadata should provide better search accuracy because they provide more semantic meaning and context, making it easier for the search engine to retrieve relevant information. In this sub-section, we have outlined a step-by-step process for preparing and indexing documents to improve search accuracy. The focus is on comparing unformatted documents to formatted ones using Quarto Markdown, and how sitemap integration enhances search engine results.

- Step 1: Document Preparation
  - Create **unformatted text**, **PDF**, or **Markdown files**.
  - Create **formatted documents** using **Quarto**, which include metadata, clear headings, and semantic structure.
- Step 2: Sitemap Generation
  - For the formatted documents, generate a **sitemap** in XML format.
  - The sitemap should list all document URLs along with relevant metadata (e.g., last modified date, frequency of changes).
- Step 3: Set Up Search Engine
  - Choose a simple search engine library, such as **Whoosh**.
  - Create a **search index** for both sets of documents (formatted and unformatted).
  - Ensure that metadata is included in the search index for the formatted documents.
- Step 4: Develop Web Crawler
  - Write a simple **web crawler** to crawl both unformatted and formatted documents.
  - For the formatted documents, ensure the crawler uses the **sitemap** to guide the indexing process.
- Step 5: Test Search Accuracy
  - Perform search queries for common terms in both unformatted and formatted datasets.
  - Measure the relevance of search results using metrics like **precision**, **recall**, and **F1 score**.
- Step 6: Analyze Results
  - Compare the performance of the search engine on unformatted versus formatted documents.
  - **Hypothesis:** Documents with structure and a sitemap will produce better search accuracy, yielding higher relevance in the results.

This Quarto Markdown setup can be used in a Jupyter notebook (ipynb) under a single section, maintaining clarity and structure in both the notebook and final rendered outputs (e.g., HTML, PDF, or DOCX).

## 5. Conclusion

The **European Environment Agency (EEA)** recognizes the growing need for generative chatbots and natural language analysis tools to facilitate easy access to CLMS data. In response, the EEA is undertaking preparatory efforts to establish the necessary standards and infrastructure for successful chatbot integration. These activities focus on ensuring that CLMS products are **findable** and **discoverable**, enabling users, regardless of technical expertise, to access environmental data seamlessly.

A key part of this strategy is making CLMS documentation and data accessible to third-party generative AI platforms. By implementing standards for formatting and exposing information—particularly through **Quarto Markdown** and **sitemaps**—CLMS ensures that high-quality, structured data is available to chatbots and AI systems. This not only enhances product discoverability but also improves user experience, allowing chatbots to guide users through complex datasets and environmental resources.

The collaboration between CLMS and the EEA lays the groundwork for a future where AI systems can efficiently retrieve and process environmental data, supporting informed decision-making and increasing public engagement with CLMS products.

## 6. References

- [1] E. Project, “CLMS - copernicus land monitoring service.” 2024. Available: <https://land.copernicus.eu/en>
- [2] M. A. Khder, “Web scraping or web crawling: State of art, techniques, approaches and application.” *International Journal of Advances in Soft Computing & Its Applications*, vol. 13, no. 3, 2021.
- [3] B. Massimino, “Accessing online data: Web-crawling and information-scraping techniques to automate the assembly of research data,” *Journal of Business Logistics*, vol. 37, no. 1, pp. 34–42, 2016.
- [4] M. A. Kausar, V. Dhaka, and S. K. Singh, “Web crawler: A review,” *International Journal of Computer Applications*, vol. 63, no. 2, pp. 31–36, 2013.
- [5] C. Saini and V. Arora, “Information retrieval in web crawling: A survey,” in *2016 international conference on advances in computing, communications and informatics (ICACCI)*, IEEE, 2016, pp. 2635–2643.
- [6] I. Hernández, C. R. Rivero, and D. Ruiz, “Deep web crawling: A survey,” *World Wide Web*, vol. 22, pp. 1577–1610, 2019.



- [7] S. Deshmukh and K. Vishwakarma, “A survey on crawlers used in developing search engine,” in *2021 5th international conference on intelligent computing and control systems (ICICCS)*, IEEE, 2021, pp. 1446–1452.
- [8] Octoparse, “Web crawl.” 2024. Available: <https://www.octoparse.com/>
- [9] J. J. Cook, “An introduction to quarto: A versatile open-source tool for data reporting and visualization.”
- [10] S. Mati, I. Civrir, and S. I. Abba, “EviewsR: An r package for dynamic and reproducible research using EViews, r, r markdown and quarto.” *R Journal*, vol. 15, no. 2, 2023.
- [11] C. Paciorek, “An example quarto markdown file,” 2023.
- [12] I. Miroshnychenko, “QUARTO: REVOLUTIONIZING CONTENT CREATION,” *Volume editor: Vitaliy Snytyuk, Dr. Sc., Prof. Program Committee: Aldrich Chris, Andreas Pester, Frederic Mallet, Hiroshi Tanaka, Iurii Krak, Yulia Khlevna, Karsten Henke, Oleg Chertov, Oleksandr Kuchanskyi, Oleksandr Marchenko, Sándor Bozóki, Vitaliy Tsyganok, Vladimir Vovk Organizing Committee: Anatoly Anisimov, Vitaliy Snytyuk, Oleksii Bychkov, Oleh Ilarionov, Yuri*, p. 189, 2023.
- [13] R. F. Hassan and S. Hussain, “Improving the web indexing quality through a website-search engine coactions,” *International Journal of Computer and Information Technology*, vol. 3, no. 2, 2014.
- [14] M. Coe, “Website indexing,” *The Indexer: The International Journal of Indexing*, vol. 34, no. 1, pp. 20–25, 2016.