# Comprehensive Technical Guidance Report

## ContentDM API Endpoints and Parameters, Streamlit Implementation, AI Model Integration, and JSON Metadata Standards

## Executive Summary

This report provides comprehensive guidance on four critical technical areas: ContentDM API integration, Streamlit iframe implementation for URL monitoring, local AI model integration (BLIP and spaCy), and JSON data package standards for metadata organization. The research synthesizes official documentation, community best practices, and proven implementation patterns to deliver actionable guidance for developers and system architects.

## 1. ContentDM API Endpoints and Parameters

### Overview

The ContentDM API is a read-only web service provided by OCLC that enables programmatic access to digital collections and metadata in JSON or XML formats. It consists of two main API sets: Server API Functions (dmwebservices) and Website API Functions (utils).

### Base URL Structure

- **Server API**: `https://server12345.contentdm.oclc.org/digital/bl/dmwebservices/index.php`
- **Website API**: `https://cdm12345.contentdm.oclc.org/utils/`

### Key Server API Functions (dmwebservices)

**Application-Level Functions**

- **wsAPIDescribe**: Returns API version information
- Signature: `?q=wsAPIDescribe/format`
- Parameters: `format` (xml|json)

**Server-Level Functions**

- **dmGetCollectionList**: Retrieves all collections
- Signature: `?q=dmGetCollectionList/format`
- Returns: Collection aliases and metadata
- Example: `dmwebservices/index.php?q=dmGetCollectionList/json`

- **dmGetDublinCoreFieldInfo**: Returns Dublin Core field information

- **dmGetLocale**: Retrieves server locale settings
- **dmGetStopWords**: Lists search stop words

## Collection-Level Functions

- **dmGetCollectionFieldInfo**: Returns metadata field information for a collection
- Signature: `?q=dmGetCollectionFieldInfo/alias/format`

- Essential for understanding available metadata fields

- **dmGetCollectionFieldVocabulary**: Retrieves controlled vocabulary for fields

- Signature: `?q=dmGetCollectionFieldVocabulary/alias/nickname/forcedict/forcefullvoc/format`

- **dmQuery**: Performs search queries on collections

- Signature: `?q=dmQuery/alias/start/field/searchterm/sort/maxrecs/suppress/pagenum/pageptr/format`
- Maximum 1024 records per call (pagination required for larger datasets)

## Item-Level Functions

- **dmGetItemInfo**: Retrieves complete item metadata
- Signature: `?q=dmGetItemInfo/alias/pointer/format`
- Returns descriptive and administrative metadata

# Website API Functions (utils)

These functions return binary data rather than text:

- **GetFile**: Downloads specific files
- Path: `/utils/getfile/collection/alias/id/pointer/filename/name`

- **GetImage**: Retrieves scaled/cropped images

- Path: `/utils/ajaxhelper/?CISOROOT=alias&CISOPTR=pointer&action=2&DMSCALE=scale...`

- **GetThumbnail**: Fetches thumbnail images

- Path: `/utils/getthumbnail/collection/alias/id/pointer`

- **GetStream**: Handles streaming media

- Path: `/utils/getstream/collection/alias/id/pointer`

# IIIF Integration

ContentDM fully supports IIIF APIs for image-based collections:
- Base path: `/digital/iiif/`
- Example: `https://i.ytimg.com/vi/_IclA1Vg4vU/hq720.jpg?sqp=-oaymwEhCK4FEIIDSFryq4qpAxMIARU-AAAAAGAElAADIQj0AgKJD&rs=AOn4CLBuBWBA08oXedPj2gFfGl1Q64VObA

# Implementation Best Practices

1. Use HTTPS for all requests
2. Handle pagination for large datasets (1024 record limit)
3. Cache frequently accessed data
4. Implement error handling for unavailable resources
5. Test with specific collection aliases before deployment

# 2. Streamlit Iframe Implementation Techniques for URL Monitoring

## Core Implementation Methods

### Basic Iframe Embedding

```python
import streamlit.components.v1 as components

# Recommended approach
components.iframe(src="https://example.com", width=800, height=600, scrolling=True)

# Alternative with st.markdown (less secure)
st.markdown(f'<iframe src="{url}" width="100%" height="500"></iframe>',
            unsafe_allow_html=True)
```

### Dynamic URL Updates

```python
import streamlit as st
import streamlit.components.v1 as components

# Store URL in session state to prevent unwanted refreshes
if 'iframe_url' not in st.session_state:
    st.session_state.iframe_url = "https://initial.url"

# Update URL based on user input
selected_url = st.selectbox("Choose URL", url_options)
if st.button("Update"):
    st.session_state.iframe_url = selected_url

components.iframe(st.session_state.iframe_url, height=500)
```

## Auto-Refresh Implementation

### Using Streamlit Autorefresh Component

```python
from streamlit_autorefresh import st_autorefresh
import streamlit.components.v1 as components

# Auto-refresh every 5 minutes with 100 rerun limit
st_autorefresh(interval=300000, limit=100, key="autorefresh_key")

# Iframe refreshes when app reruns
components.iframe("https://dynamic-content-site.com", height=500)
```

**Manual Refresh Control**

```python
import streamlit as st

# Manual refresh button
if st.button("🔄 Refresh Content"):
    st.rerun()

# Cache-busting technique
import time
cache_buster = int(time.time()) if st.button("Force Refresh") else ""
url_with_cache_buster = f"https://example.com?_cb={cache_buster}"
```

## URL Monitoring and Tracking

### Query Parameter Management

```python
import urllib.parse

def clean_url(url):
    """Remove query parameters to prevent unintended reloads"""
    parsed = urllib.parse.urlparse(url)
    return urllib.parse.urlunparse((
        parsed.scheme, parsed.netloc, parsed.path, '', '', ''
    ))

# Use cleaned URL for stable iframe embedding
clean_iframe_url = clean_url(original_url)
components.iframe(clean_iframe_url, height=500)
```

### Status Monitoring

```python
import requests
import streamlit as st

def check_url_status(url):
    try:
        response = requests.head(url, timeout=5)
        return response.status_code == 200
    except:
        return False

url = st.text_input("Enter URL to monitor")
if url:
    status = "🟢 Online" if check_url_status(url) else "🔴 Offline"
    st.write(f"Status: {status}")

    if check_url_status(url):
        components.iframe(url, height=500)
```

## Best Practices for Iframe Implementation

1. **Security**: Use `components.iframe()` over raw HTML for better security
2. **Performance**: Implement caching with `@st.cache_data` for expensive operations
3. **Responsiveness**: Use percentage-based widths ( `width="100%"` ) for responsive design
4. **Accessibility**: Set appropriate `scrolling` and `tab_index` parameters
5. **Error Handling**: Always check URL accessibility before embedding

6. **Session Management**: Use `st.session_state` to persist iframe states across reruns

## Embedding Streamlit Apps

For embedding Streamlit apps themselves:

```html
<iframe src="https://yourapp.streamlit.app?embed=true"
        style="height: 450px; width: 100%;"></iframe>
```

Use embed options for customization:
- `?embed=true&embed_options=show_toolbar`
- `?embed=true&embed_options=disable_scrolling`

---

# 3. Local AI Models Integration: BLIP and spaCy in Streamlit

## BLIP Image Captioning Integration

### Environment Setup

```
pip install streamlit transformers torch pillow
```

### Basic BLIP Implementation

```python
import streamlit as st
from PIL import Image
from transformers import BlipProcessor, BlipForConditionalGeneration

@st.cache_resource
def load_blip_model():
    """Load BLIP model with caching for performance"""
    processor = BlipProcessor.from_pretrained("Salesforce/blip-image-captioning-base")
    model = BlipForConditionalGeneration.from_pretrained("Salesforce/blip-image-cap-
tioning-base")
    return processor, model

def main():
    st.title("BLIP Image Captioning App")

    processor, model = load_blip_model()

    uploaded_file = st.file_uploader("Upload an image", type=["jpg", "jpeg", "png"])

    if uploaded_file:
        image = Image.open(uploaded_file).convert('RGB')
        st.image(image, caption="Uploaded Image", use_column_width=True)

        if st.button("Generate Caption"):
            with st.spinner("Generating caption..."):
                inputs = processor(image, return_tensors="pt")
                out = model.generate(**inputs, max_length=50, num_beams=5)
                caption = processor.decode(out[0], skip_special_tokens=True)
                st.success(f"Generated Caption: {caption}")
```

**GPU/CPU Optimization**

```python
import torch

@st.cache_resource
def load_optimized_blip_model():
    device = "cuda" if torch.cuda.is_available() else "cpu"
    processor = BlipProcessor.from_pretrained("Salesforce/blip-image-captioning-base")
    model = BlipForConditionalGeneration.from_pretrained("Salesforce/blip-image-cap-
tioning-base")
    model = model.to(device)

    st.info(f"Model loaded on: {device}")
    return processor, model, device

# Usage with device awareness
processor, model, device = load_optimized_blip_model()
inputs = processor(image, return_tensors="pt").to(device)
```

# spaCy NER Integration

**Installation and Setup**

```
pip install streamlit spacy spacy-streamlit
python -m spacy download en_core_web_sm
```

**Basic NER Implementation**

```python
import streamlit as st
import spacy
from spacy_streamlit import visualize_ner
import pandas as pd

@st.cache_resource
def load_spacy_model():
    """Load spaCy model with caching"""
    return spacy.load("en_core_web_sm")

def main():
    st.title("Named Entity Recognition with spaCy")

    nlp = load_spacy_model()

    # Model selection
    model_options = ["en_core_web_sm", "en_core_web_md", "en_core_web_lg"]
    selected_model = st.sidebar.selectbox("Select Model", model_options)

    # Text input
    text = st.text_area(
        "Enter text for NER analysis",
        "Apple Inc. is headquartered in Cupertino, California. Tim Cook is the CEO."
    )

    if st.button("Analyze Text"):
        doc = nlp(text)

        # Visualization using spacy-streamlit
        st.subheader("Entity Visualization")
        visualize_ner(doc, labels=nlp.get_pipe("ner").labels)

        # Custom entity table
        st.subheader("Extracted Entities")
        entities = [(ent.text, ent.start_char, ent.end_char, ent.label_)
                    for ent in doc.ents]

        if entities:
            df = pd.DataFrame(entities, columns=["Text", "Start", "End", "Label"])
            st.dataframe(df)
        else:
            st.info("No entities found in the text.")
```

**Advanced NER with URL Processing**

```python
from bs4 import BeautifulSoup
import requests

def extract_text_from_url(url):
    """Extract text content from URL for NER analysis"""
    try:
        response = requests.get(url, timeout=10)
        soup = BeautifulSoup(response.content, 'html.parser')

        # Extract text from paragraphs
        paragraphs = soup.find_all('p')
        text = ' '.join([p.get_text() for p in paragraphs])
        return text[:5000]  # Limit text length
    except Exception as e:
        st.error(f"Error extracting text from URL: {e}")
        return None

# Integration in Streamlit app
url = st.text_input("Enter URL for NER analysis")
if url and st.button("Analyze URL"):
    text = extract_text_from_url(url)
    if text:
        doc = nlp(text)
        visualize_ner(doc, labels=nlp.get_pipe("ner").labels)
```

# Performance Optimization Strategies

## Model Caching and Memory Management

```python
import gc

@st.cache_resource
def get_model_pipeline():
    """Cached model loading with memory optimization"""
    nlp = spacy.load("en_core_web_sm", exclude=["parser", "tagger"])  # Load only NER
    return nlp

def clear_memory():
    """Clear memory after processing large batches"""
    gc.collect()
    if torch.cuda.is_available():
        torch.cuda.empty_cache()
```

**Batch Processing**

```python
def process_multiple_texts(texts, nlp_model):
    """Efficient batch processing for multiple texts"""
    docs = list(nlp_model.pipe(texts, disable=["parser", "tagger"]))
    return docs

# Usage in Streamlit
uploaded_files = st.file_uploader("Upload text files", accept_multiple_files=True)
if uploaded_files and st.button("Process All Files"):
    texts = [file.read().decode() for file in uploaded_files]
    docs = process_multiple_texts(texts, nlp)

    for i, doc in enumerate(docs):
        st.subheader(f"File {i+1} Results")
        visualize_ner(doc)
```

## Concurrency Considerations

**Thread-Safe Model Loading**

```python
import threading

class ModelManager:
    _instance = None
    _lock = threading.Lock()

    def __new__(cls):
        if cls._instance is None:
            with cls._lock:
                if cls._instance is None:
                    cls._instance = super().__new__(cls)
                    cls._instance.models = {}
        return cls._instance

    def get_model(self, model_name):
        if model_name not in self.models:
            if model_name == "blip":
                processor = BlipProcessor.from_pretrained("Salesforce/blip-image-cap-
tioning-base")
                model = BlipForConditionalGeneration.from_pretrained("Salesforce/blip-
image-captioning-base")
                self.models[model_name] = (processor, model)
            elif model_name == "spacy":
                self.models[model_name] = spacy.load("en_core_web_sm")

        return self.models[model_name]

# Usage
model_manager = ModelManager()
nlp = model_manager.get_model("spacy")
```

# 4. JSON Data Package Standards for Metadata Organization

## Frictionless Data Package Standards

### Basic Data Package Structure

Frictionless Data uses JSON extensively for metadata organization through the Data Package specification:

```json
{
  "name": "my-dataset",
  "title": "My Dataset",
  "description": "A comprehensive dataset for analysis",
  "version": "1.0.0",
  "keywords": ["dataset", "analysis", "research"],
  "licenses": [
    {
      "name": "CC-BY-4.0",
      "path": "https://creativecommons.org/licenses/by/4.0/",
      "title": "Creative Commons Attribution 4.0"
    }
  ],
  "resources": [
    {
      "name": "data",
      "path": "data.csv",
      "title": "Main Dataset",
      "description": "Primary data file",
      "format": "csv",
      "mediatype": "text/csv",
      "schema": {
        "fields": [
          {
            "name": "id",
            "type": "integer",
            "description": "Unique identifier"
          },
          {
            "name": "name",
            "type": "string",
            "description": "Record name"
          }
        ]
      }
    }
  ]
}
```

### Key Metadata Properties

- **Required**: `name`, `resources`
- **Recommended**: `id`, `title`, `description`, `version`, `keywords`, `licenses`
- **Optional**: `contributors`, `created`, `homepage`, `image`

### JSON Resource Handling

For JSON data resources within packages:

```json
{
  "name": "json-resource",
  "path": "data.json",
  "format": "json",
  "encoding": "utf-8",
  "profile": "tabular-data-resource",
  "schema": {
    "fields": [
      {
        "name": "timestamp",
        "type": "datetime",
        "format": "iso"
      }
    ]
  }
}
```

## JSON-LD Best Practices

### Context and Vocabulary Usage

```json
{
  "@context": "https://schema.org/",
  "@type": "Dataset",
  "@id": "https://example.org/dataset/123",
  "name": "Research Dataset",
  "description": "Comprehensive research data",
  "creator": {
    "@type": "Person",
    "name": "Jane Researcher",
    "@id": "https://orcid.org/0000-0000-0000-0000"
  },
  "datePublished": "2023-01-15",
  "keywords": ["research", "data", "analysis"],
  "license": "https://creativecommons.org/licenses/by/4.0/",
  "spatialCoverage": {
    "@type": "Place",
    "name": "Global"
  },
  "temporalCoverage": "2020/2023"
}
```

### Entity Identification and Typing

Best practices for JSON-LD metadata organization:

1. **Use Unique Identifiers**: Always include `@id` for entities
2. **Specify Types**: Include `@type` for all objects
3. **Reference External Entities**: Use URIs for people, places, organizations
4. **Structured Properties**: Organize related data in nested objects
5. **Standardized Vocabularies**: Prefer Schema.org, Dublin Core terms

**Multi-Schema Integration**

```json
{
  "@context": [
    "https://schema.org/",
    {
      "dcterms": "http://purl.org/dc/terms/",
      "foaf": "http://xmlns.com/foaf/0.1/"
    }
  ],
  "@type": "Dataset",
  "name": "Integrated Metadata Example",
  "dcterms:creator": {
    "@type": "foaf:Person",
    "foaf:name": "Research Team",
    "foaf:mbox": "team@example.org"
  },
  "dcterms:subject": ["metadata", "standards", "integration"],
  "schema:dateCreated": "2023-01-01"
}
```

## Dublin Core in JSON Format

### Basic Dublin Core JSON-LD Implementation

```json
{
  "@context": {
    "dc": "http://purl.org/dc/terms/",
    "dcmitype": "http://purl.org/dc/dcmitype/"
  },
  "@type": "dcmitype:Dataset",
  "dc:title": "Dublin Core Metadata Example",
  "dc:creator": [
    {
      "@type": "http://purl.org/dc/terms/Agent",
      "dc:name": "Primary Author"
    }
  ],
  "dc:subject": ["metadata", "dublin core", "standards"],
  "dc:description": "An example of Dublin Core metadata in JSON-LD format",
  "dc:publisher": "Example Organization",
  "dc:date": "2023-01-01",
  "dc:type": "Dataset",
  "dc:format": "application/json",
  "dc:identifier": "https://doi.org/10.1000/example",
  "dc:language": "en",
  "dc:rights": "CC BY 4.0",
  "dc:coverage": {
    "spatial": "Global",
    "temporal": "2020-2023"
  }
}
```

### Dublin Core Elements in JSON

The 15 core Dublin Core elements in JSON structure:

```json
{
  "@context": "http://purl.org/dc/terms/",
  "title": "Resource Title",
  "creator": "Content Creator",
  "subject": ["keyword1", "keyword2"],
  "description": "Resource description",
  "publisher": "Publishing Organization",
  "contributor": "Additional Contributors",
  "date": "2023-01-01",
  "type": "Text",
  "format": "application/pdf",
  "identifier": "ISBN:123456789",
  "source": "Source Publication",
  "language": "en",
  "relation": "Related Resource URI",
  "coverage": "Geographic or Temporal Scope",
  "rights": "Rights Information"
}
```

## Metadata Organization Patterns

### Hierarchical Resource Description

```json
{
  "name": "research-project",
  "title": "Comprehensive Research Project",
  "resources": [
    {
      "name": "raw-data",
      "title": "Raw Data Files",
      "path": "data/raw/",
      "resources": [
        {
          "name": "survey-responses",
          "path": "data/raw/survey.csv",
          "schema": "schemas/survey-schema.json"
        }
      ]
    },
    {
      "name": "processed-data",
      "title": "Processed Data",
      "path": "data/processed/",
      "derivedFrom": "raw-data"
    }
  ]
}
```

**Provenance and Lineage Tracking**

```json
{
  "@context": {
    "prov": "http://www.w3.org/ns/prov#",
    "schema": "https://schema.org/"
  },
  "@type": "schema:Dataset",
  "name": "Processed Dataset",
  "prov:wasDerivedFrom": {
    "@id": "https://example.org/raw-dataset",
    "@type": "schema:Dataset"
  },
  "prov:wasGeneratedBy": {
    "@type": "prov:Activity",
    "prov:startedAtTime": "2023-01-01T10:00:00Z",
    "prov:used": [
      "https://example.org/processing-script",
      "https://example.org/configuration"
    ]
  }
}
```

# Key Takeaways and Best Practices

### ContentDM API

- **Always use HTTPS** for API requests to OCLC-hosted instances
- **Implement pagination** for queries returning more than 1024 records
- **Cache frequently accessed data** to improve performance
- **Use collection aliases** correctly - they're essential for all collection-level operations
- **Handle errors gracefully** - implement fallbacks for unavailable resources

### Streamlit Iframe Implementation

- **Prefer** `st.components.v1.iframe()` over raw HTML for better security
- **Use session state** to persist iframe URLs and prevent unwanted refreshes
- **Implement auto-refresh** with the `streamlit-autorefresh` component for dynamic content
- **Clean URLs** to remove unnecessary query parameters that cause reloads
- **Add status monitoring** to check URL accessibility before embedding

### AI Model Integration

- **Use** `@st.cache_resource` for model loading to prevent repeated initialization
- **Implement GPU/CPU detection** with appropriate fallbacks
- **Optimize memory usage** by excluding unused model components
- **Handle batch processing** efficiently for multiple inputs
- **Provide clear user feedback** during model loading and inference

### JSON Metadata Standards

- **Use standardized vocabularies** like Schema.org and Dublin Core terms
- **Provide unique identifiers** (@id) for all entities in JSON-LD
- **Implement proper typing** with @type declarations

- **Structure data hierarchically** for complex resources
- **Validate metadata** using appropriate tools and schemas

## General Development Practices

- **Cache expensive operations** using Streamlit's caching decorators
- **Implement error handling** with try-catch blocks and user-friendly messages
- **Use environment variables** for configuration and API keys
- **Test across different environments** before deployment
- **Document API integrations** thoroughly for maintenance

---

# References

## ContentDM API Documentation

1. OCLC ContentDM API Reference - https://help.oclc.org/Metadata_Services/CONTENTdm/Advanced_website_customization/API_Reference/CONTENTdm_API
2. ContentDM Server API Functions - https://help.oclc.org/Metadata_Services/CONTENTdm/Advanced_website_customization/API_Reference/CONTENTdm_API/CONTENTdm_Server_API_Functions_dmwebservices
3. ContentDM Website API Reference - https://help.oclc.org/Metadata_Services/CONTENTdm/Advanced_website_customization/API_Reference/CONTENTdm_API/CONTENTdm_Website_API_Reference_utils
4. ContentDM Tips and Tricks - https://evanwill.github.io/_drafts/notes/contentdm-tips.html
5. UWM Libraries ContentDM API Guide - https://guides.library.uwm.edu/c.php?g=1066592&p=7802347

## Streamlit Implementation

1. Streamlit Components iframe Documentation - https://docs.streamlit.io/develop/api-reference/custom-components/st.components.v1.iframe
2. Streamlit App Embedding Guide - https://docs.streamlit.io/deploy/streamlit-community-cloud/share-your-app/embed-your-app
3. Streamlit Community Discussions on iframe - https://discuss.streamlit.io/t/auto-refresh-an-iframe-inside-streamlit/42168
4. Streamlit Autorefresh Component - https://discuss.streamlit.io/t/streamlit-autorefresh/14519

## AI Model Integration

1. BLIP Model Integration Tutorial - https://www.marktechpost.com/2025/03/13/a-coding-guide-to-build-a-multimodal-image-captioning-app-using-salesforce-blip-model-streamlit-ngrok-and-hugging-face/
2. Hugging Face BLIP Documentation - https://huggingface.co/Salesforce/blip-image-captioning-base
3. spaCy-Streamlit GitHub Repository - https://github.com/explosion/spacy-streamlit
4. spaCy NER with Streamlit Tutorial - https://towardsdatascience.com/build-a-named-entity-recognition-app-with-streamlit-f157672f867f
5. GPU vs CPU in Streamlit - https://docs.vultr.com/how-to-deploy-a-deep-learning-model-with-streamlit

## JSON Metadata Standards

1. Frictionless Data Package Specification - https://specs.frictionlessdata.io/data-package/

2. JSON-LD Best Practices - https://w3c.github.io/json-ld-bp/

3. Schema.org JSON-LD Documentation - https://developers.google.com/search/docs/appearance/structured-data/intro-structured-data

4. Dublin Core Metadata Terms - https://www.dublincore.org/specifications/dublin-core/dcmi-terms/

5. Dublin Core JSON-LD Examples - https://www.dublincore.org/resources/metadata-basics/