

Trustworthy AI-based Multi-Modal Data Processing for Cultural Heritage Monitoring and Preservation

Lingxiao Kong¹, Athos Agapiou², George Pavlidis³, Zeyd Boukhers¹

¹Fraunhofer Institute for Applied Information Technology FIT, ²Cyprus University of Technology, ³ATHENA Research Center

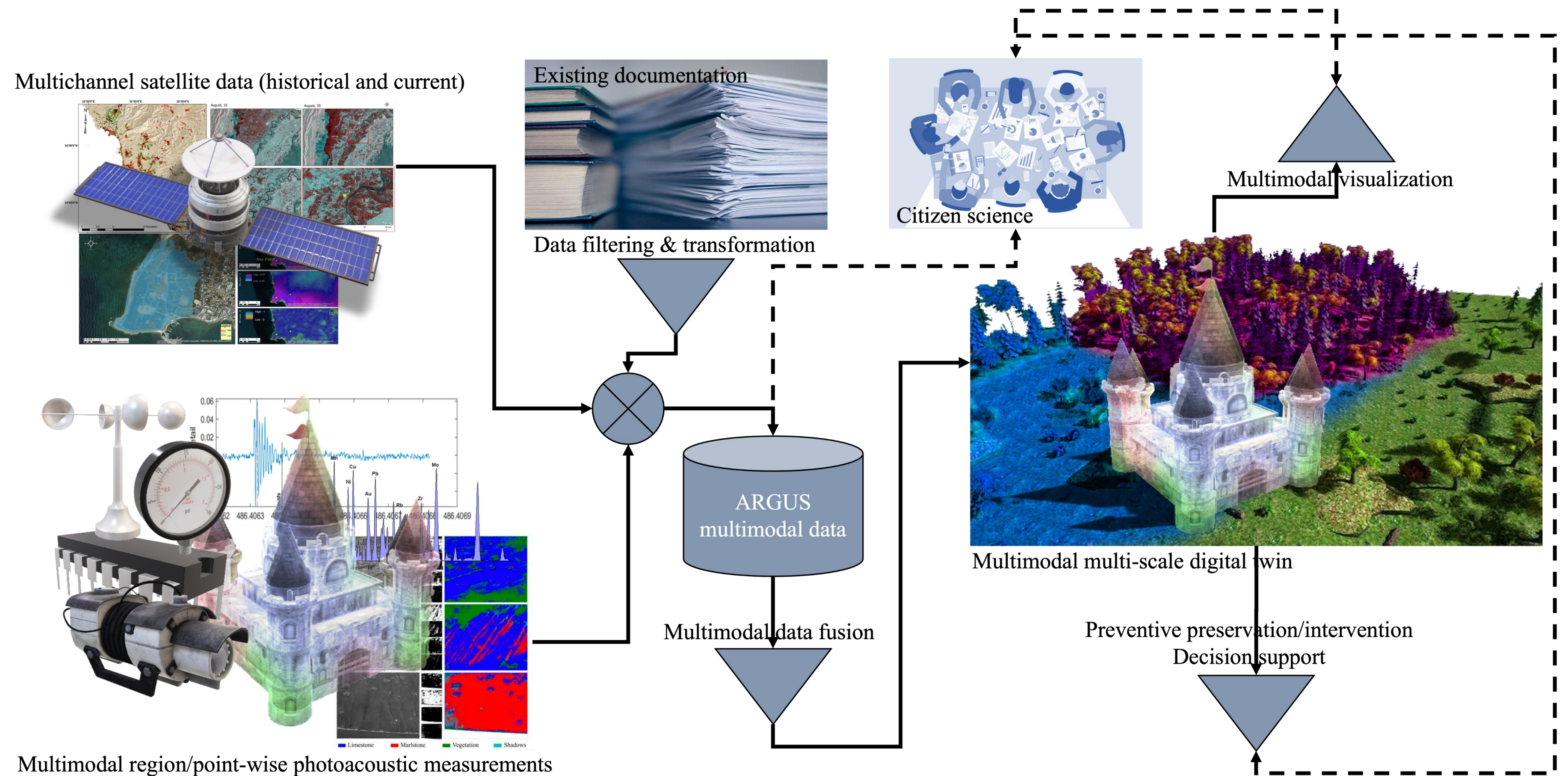
1. CH Data Challenges

Heterogeneous data of diverse types and modalities are continuously fed into the ARGUS digital twin, requiring sophisticated processing before they can be effectively utilized. Key challenges include:

- Inconsistent formats and non-interoperable standards
- Missing or ambiguous metadata
- Limited spatial and temporal coverage
- Incompatibility for multi-modal AI analysis
- Misalignment with FAIR principles

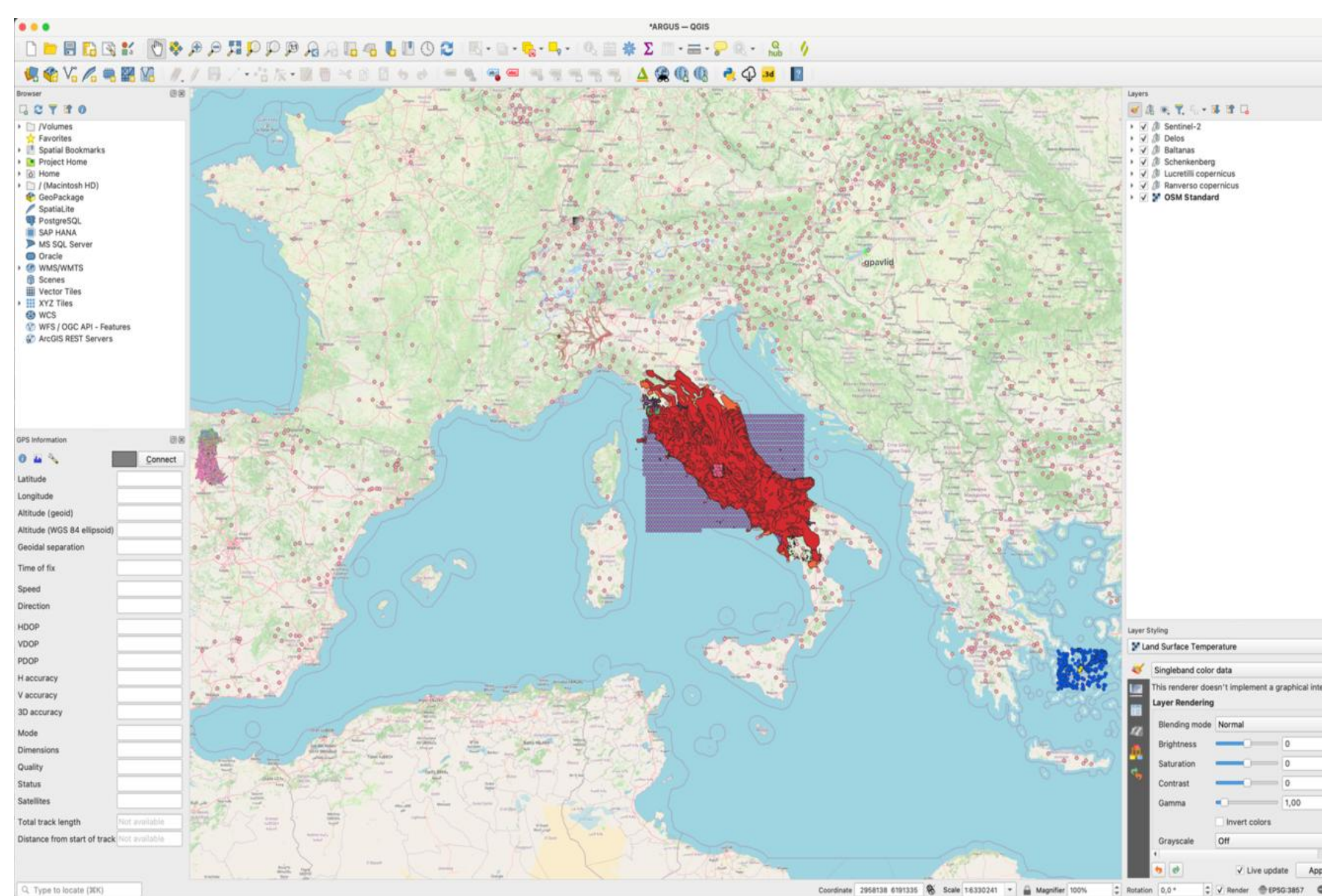
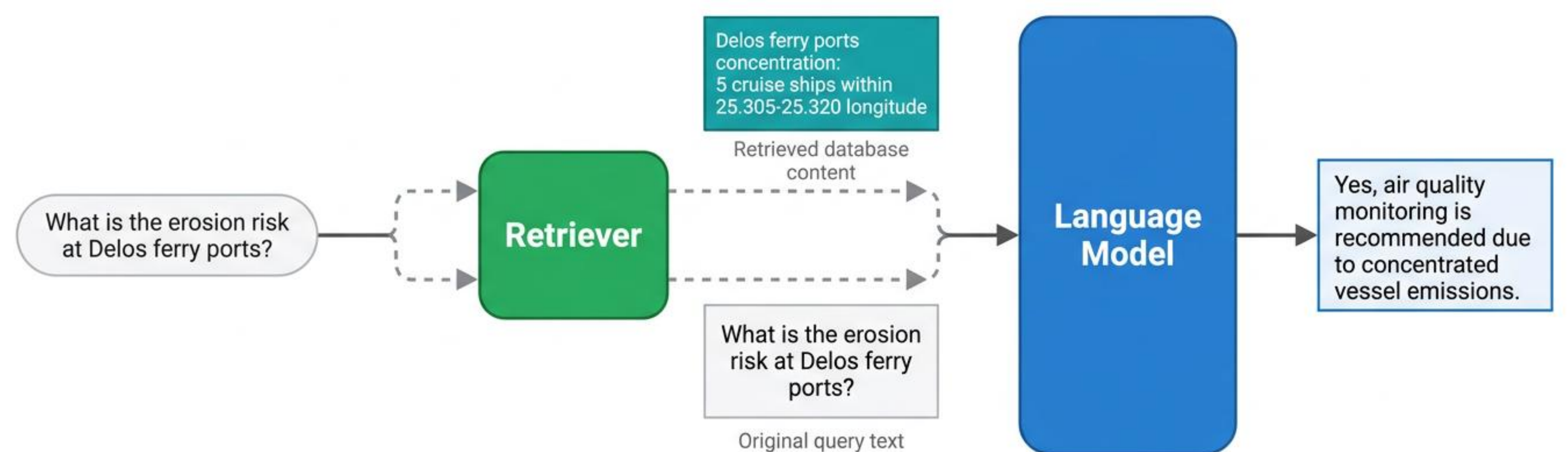
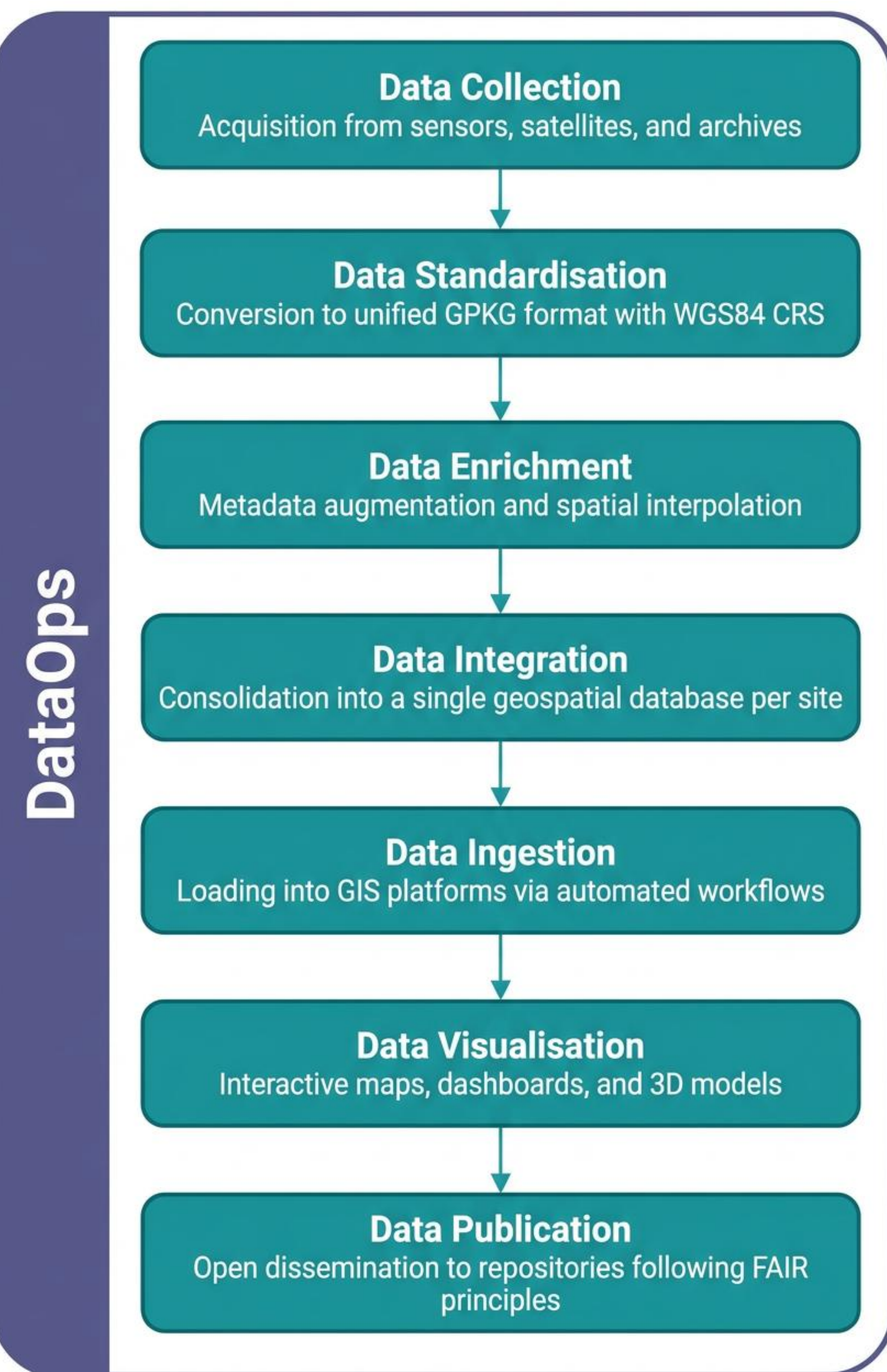
2. Data Processing Pipeline

The CH data processing pipeline comprising **seven sequential stages** from data sources through AI-ready data, governed by an overarching DataOps framework. Each stage progressively transforms raw heterogeneous data spanning geometric, visual, and documentary modalities into standardized, enriched, and queryable resources, with DataOps ensuring automation, monitoring, and continuous quality control throughout the entire lifecycle, providing a reliable data foundation for AI-driven heritage monitoring and preservation.



3. RAG-based LLM Application

The application supports two functionalities. For **database querying**, natural language queries are fed into a Retriever that retrieves relevant records from the internal GPKG database, which are jointly provided to the Language Model to generate grounded, contextually accurate responses. For **attribute annotation**, the same architecture retrieves definitions and unit references from external knowledge bases such as WikiData to generate missing attribute descriptions. All outputs are labelled [GENERATED] with source citations to ensure transparency and traceability.



4. Processing Results and Further Steps

QGIS visualization of the integrated ARGUS geospatial database and key processing results, demonstrating that the proposed framework effectively transforms heterogeneous cultural heritage data into a standardized, enriched, and queryable resource, providing a reliable foundation for **downstream AI solutions** including anomaly detection, visualization, and decision support for heritage monitoring and preservation.

DATASETS INTEGRATED 249 → 1 ADF, CSV, GDB, GPKG, SHP, TIF, XLS, XLSX	PIPELINE SUCCESS RATE 98.8% 246 of 249 succeeded
PILOT SITES COVERED 5 Schenkenberg, Lucrettili, Delos, Ranverso, Baltanás	OUTPUT FORMAT GPKG/WGS84 CRS
TEMPORAL COVERAGE days → sec Nearest-neighbour imputation	SPATIAL COVERAGE 22% → 76% Nearest-neighbour imputation
ENRICHED METADATA 14% → 100% Descriptions and units added	QUERY EFFICIENCY hours → minutes LLM-supported querying

