

WP2 – Data-Lake Infrastructure

Draft Review

MAPS | Multidimensional Area Place-based System

Tech team bi-weekly · March 2, 2026

Agenda

What we're presenting

1. WP2 scope and deliverables
2. Medallion architecture (Bronze → Silver → Gold)
3. Technology stack
4. Cloud architecture
5. Specific technical challenges
6. Current draft status

What we want to discuss

- Open architectural decisions
- Deployment: managed vs self-hosted
- MVP Phase 1a priorities
- Feedback on the document

WP2 – Scope and Deliverables

Task 2.1 · M9–M10

Technical architecture

- ✓ D2.1.1 – Data-Lake Technical Design
- ✓ D2.1.2 – Cloud Solution Design
- ✓ D2.1.3 – Hosting Specifications

Deadline: 31/03/2026

Task 2.2 · M10–M12

ETL pipeline

- ○ D2.2 – ETL scripts (Python/SQL)

Deadline: 31/05/2026

Task 2.3 · M10–M12

Validation

- ○ D2.3 – ETL Validation Report

Deadline: 31/05/2026

Today we present the draft of the three Task 2.1 deliverables — architecture, cloud, hosting. D2.2 and D2.3 will be completed during the development phase.

D2.1.1 – Requirements

FUNCTIONAL

RF1 – Heterogeneous acquisition

180+ datasets from ISTAT, Ministries, public registries. Formats: CSV, XLSX, PDF, HTML.

RF2 – Historical series management

2010–2025, ~8,000 municipalities. Must handle COVID gaps and ~100 municipal mergers/splits.

RF3 – Standardisation

ISTAT code normalisation, spatial reference EPSG:32632, distinction between missing values and semantic nulls.

RF4 – Multi-layer architecture

Progressive data refinement: raw → clean → application-ready.

NON-FUNCTIONAL

RNF1 – Completeness $\geq 95\%$ of municipalities per priority dataset

RNF2 – Accuracy $\geq 99.9\%$ of records with valid ISTAT code

RNF3 – Timeliness Priority pipelines complete within 24h

RNF4 – Lineage Full traceability for every record, source → Bronze → Silver → Gold

RNF5 – Scalability Up to 1TB over 3 years, ≥ 100 records/sec throughput

RNF6 – Governance Full dataset cataloging, real-time quality dashboard, audit logs

D2.1.1 – Architectural Constraints

V1 – Self-hosted open source

No managed cloud services (BigQuery, Azure, AWS excluded). Full stack: PostgreSQL, PostGIS, Prefect, OpenMetadata, DuckDB. Vendor independence and predictable costs.

V3 – Territorial standards

PostGIS as de facto standard for spatial data in PostgreSQL. Reference system EPSG:32632 (WGS84 / UTM zone 32N), compatible with all Italian national geographic datasets.

V2 – Scale adequacy

Data volumes $\sim 10^4$ rows \times $\sim 10^3$ columns. Enterprise-scale infrastructure (BigQuery, Spark) would be overkill — unjustified cost and complexity for actual workloads.

V4 – FAIR compliance

Released datasets must be Findable (structured metadata), Accessible (CC-BY/CC0), Interoperable (GeoJSON, GeoParquet, CSV), Reusable (DOI for scientific citability).

Given these constraints, the chosen technical solution for RF4 is the **Medallion pattern** (Bronze → Silver → Gold), detailed in deliverable D2.1.1.

Medallion Architecture

EXTERNAL SOURCES → BRONZE → SILVER → GOLD → APPLICATIONS

Bronze

Immutable archive

- Local file system
- Original files, no transformations
- Write-once, never-modify
- SHA-256 integrity checksums
- Sidecar JSON metadata files

Silver

Clean and validated data

- PostgreSQL schema `silver`
- EAV (Entity-Attribute-Value) schema
- SCD Type 2 for municipal changes
- Great Expectations validation
- Single source of truth

Gold

Application-ready data

- PostgreSQL + PostGIS
- Denormalized wide tables
- Spatial geometries ready to query
- Feeds API, dashboards, SLO algorithms
- DuckDB for ad-hoc analytics

Silver Layer – Why EAV?

The problem: 180+ datasets with heterogeneous structures — traditional schema requires hundreds of columns and constant migrations. **The solution:** each observation becomes a row (municipality, attribute, value).

Traditional schema (problematic)

codice_istat	population	nurseries	has_hospital	ato_gas	...
001001	45230	12	true	?	...
001002	8420	2	false	?	...

Every new dataset requires new columns → ALTER TABLE

EAV schema (adopted)

entity	attribute	value	valid_from	source
001001	population	45230	2024-01-01	ISTAT
001001	nurseries	12	2024-01-01	ISTAT
001001	has_hospital	true	2015-01-01	MinSalute
001002	population	8420	2024-01-01	ISTAT

New datasets = new rows. Schema stays stable.

Main trade-off: more complex queries (PIVOT with CASE WHEN), but complexity is confined to the Silver → Gold transformation layer, invisible to data mart users.

Technical Challenges

S1 – No point-level geolocation

Most services (schools, hospitals, tobacconists) lack precise coordinates. Approach: presence/absence at municipal level. Isochrones operate on municipality-to-municipality distances.

S2 – Administrative boundary changes

~100 municipal mergers/splits in 2010-2025. Solution: `comuni_variazioni` table with SCD Type 2 and `valid_from` / `valid_to` columns.

S3 – COVID temporal gaps

Methodological discontinuities in 2020-2021.

`covid_affected` flag on impacted datasets. Options: statistical imputation or explicit exclusion of affected periods.

S4 – Format heterogeneity

207 datasets across CSV, XLSX, HTML, PDF. Multi-layer extraction stack: Docling (PDF, 97.9% accuracy), Pandas (CSV/Excel), BeautifulSoup (HTML).

Technology Stack

Component	Technology	Role
Primary storage	PostgreSQL 17 + PostGIS 3.5	Master data, spatial operations
Orchestration	Prefect 3.x	Scheduling, monitoring, ETL
Analytics	DuckDB 1.x	Federated queries from PostgreSQL
Internal governance	OpenMetadata 1.x	Catalog, lineage, quality
Open data	CKAN 2.11	Public catalog (DCAT-AP_IT)
Data quality	Great Expectations 1.x	Validation, anomaly detection

Guiding principles

- **Open source self-hosted:** no vendor lock-in, predictable costs
- **Scale adequacy:** $\sim 10^4$ rows \times $\sim 10^3$ columns. BigQuery is overkill.
- **Territorial standards:** PostGIS is the de facto standard for spatial data
- **FAIR compliance:** CC-BY, GeoJSON/GeoParquet, DOI

D2.1.2 – Cloud Solution Design

Provider: DigitalOcean

- Fixed monthly pricing — no variable billing
- Managed Kubernetes (DOKS) — control plane managed by DO
- Simpler learning curve vs AWS/GCP/Azure hyperscalers

Why Kubernetes over Docker Compose

- Declarative infrastructure, self-healing
- Workload isolation per service (CPU/memory limits)
- Supports project growth without architectural migration

Managed PostgreSQL vs self-hosted

- Managed: automatic backup, failover, PITR — zero ops overhead
- Self-hosted (CNPG operator): full control, lower cost, higher ops burden
- D2.1.2 presents both — decision deferred to D2.1.3

D2.1.2 cost estimate (Managed PostgreSQL)

Item	Cost/month
DOKS control plane	\$12
3 nodes (4 vCPU / 8 GB)	\$192
Managed PostgreSQL	\$80
Block storage (150 GB)	\$15
Load Balancer	\$12
Total	~\$311/month

Current MVP deployment: op-linkurious server with Docker Compose. Migration to DOKS planned for the production phase.

D2.1.3 – Hosting Specifications

laC: Pulumi (Python)

- Full cluster provisioned with `pulumi up --stack maps`
- Two node pools: `general-purpose` (2× 2vCPU/4GB, fixed) for system components; `workloads` (1–3× 4vCPU/8GB, autoscaling) for Prefect, PostgreSQL, OpenMetadata, CKAN

Database: CloudNative PostgreSQL (CNPB)

- Self-hosted operator — resolves D2.1.2 open question
- 2 instances, synchronous replication, automatic failover
- WAL backup to S3, point-in-time recovery (30-day retention)

Secrets: AWS Parameter Store + External Secrets Operator

- Path schema: `/maps/{service}/{parameter}`
- Synced to Kubernetes Secrets every hour

D2.1.3 cost estimate (CNPB self-hosted + autoscaling)

Item	Cost/month
DOKS control plane	\$12
general-purpose pool (2 nodes)	\$48
workloads pool (min 1 node)	\$48
workloads pool (max 3 nodes)	\$144
Block storage (200 GB)	\$20
Load Balancer	\$12
Spaces backup	\$5
Total (idle)	~\$145/month
Total (ETL peak)	~\$241/month

\$70–166/month less than D2.1.2 estimate thanks to CNPB self-hosted + autoscaling.

MVP Phase 1a – Critical Datasets

5 minimum datasets to bootstrap the SLO algorithms (WP3):

Dataset	Source	Format	Frequency	Priority
Municipal boundaries + metadata	ISTAT	Shapefile/JSON	Annual	High
Resident population	ISTAT	CSV	Annual	High
Commuting matrices	ISTAT	CSV	Decennial	High
Accessibility to transport infrastructure	ISTAT	XLSX	Spot	Medium
Healthcare facilities	Min. Salute	Excel	Annual	Medium

≥95%

municipality coverage per priority dataset

≥99.9%

records with valid ISTAT code

<24h

latency for priority pipelines

Draft Status

Completed (Task 2.1)

- **D2.1.1 – Data-Lake Architecture**
 - Medallion pattern with rationale
 - Detailed Silver EAV schema
 - Gold data marts (full SQL)
 - Governance and lineage design
 - Best practices (idempotency, schema evolution)
- **D2.1.2 – Cloud Solution Design**
 - Provider selection (DigitalOcean)
 - Cluster sizing and cost estimates
 - Managed vs self-hosted PostgreSQL
- **D2.1.3 – Hosting Specifications**
 - Detailed server specifications
 - Service configuration

To be developed (Task 2.2 / 2.3)

- **D2.2 – ETL Scripts**
 - Multi-worker-pool structure defined
 - Common pipeline pattern documented
 - Implementation: M10–M12
- **D2.3 – Validation Report**
 - Acceptance criteria defined
 - Test suite structured
 - Execution: after implementation

Question: is the current draft sufficient for the formal review of deliverables D2.1.x by 31/03/2026?

Open Questions – Project & Team

DOMAIN & IDENTITY

Q1 – What is the final domain? D2.1.3 uses `maps.gransassotech.it` as placeholder.

Q2 – Official project name: "MAPS" or "GST-MAPS"? This affects the CKAN catalog and Zenodo DOI metadata.

Q3 – Is there already a dedicated DigitalOcean account for MAPS, or does one need to be created?

Q4 – AWS account for Route53: dedicated to MAPS, or use the existing DEPP/Openpolis account?

TEAM & CONTINUITY

Q5 – Who manages the infrastructure after the project ends (May 2027)? DigitalOcean costs ~\$300/month ongoing.

Q6 – Do datasets and the web platform remain publicly accessible after closure? Is there a sustainability plan?

Q7 – Who builds the frontend (scroll-telling site + data visualization)? Is a separate contract planned?

Open Questions – Technical & Process

TECHNICAL DECISIONS

Q8 – Final database: DigitalOcean managed or self-hosted via CloudNativePG? D2.1.2 presents both without a final decision.

Q9 – Is CKAN confirmed as the public open data catalog? It requires non-trivial setup and maintenance — are resources allocated?

Q10 – D2.1.3 resolves the managed vs self-hosted question in favour of CloudNative PostgreSQL (CNPG): lower cost (\$145–241/month vs \$311), full PostGIS control, higher ops burden. Does the team confirm this choice?

DELIVERABLES & FORMAL PROCESS

Q11 – Task 2.2 (ETL development) starts M10, but it requires a running infrastructure: PostgreSQL, Prefect, Bronze storage. Who sets this up, and when? There is no task in the plan explicitly covering infrastructure deployment before development begins.

Q12 – Must D2.1.x be formally submitted to MUR by 31/03/2026, or is that an internal deadline with a GST acceptance phase first?

Q13 – Is there a standard deliverable format required by the FRES 2 call? Template, cover page, signature, submission protocol?

Next Steps

By 31/03

- Internal review of D2.1.x
- Close today's feedback
- Finalize document

M10–M11 (Mar–Apr)

- Start ETL pipeline development
- 5 critical MVP datasets
- Prefect + PostgreSQL setup

M12 (May 31)

- D2.2 – ETL scripts complete
- D2.3 – Validation report
- Hand-off to WP3 (SLO)

Full draft: `deliverables/docs/wp2/`