



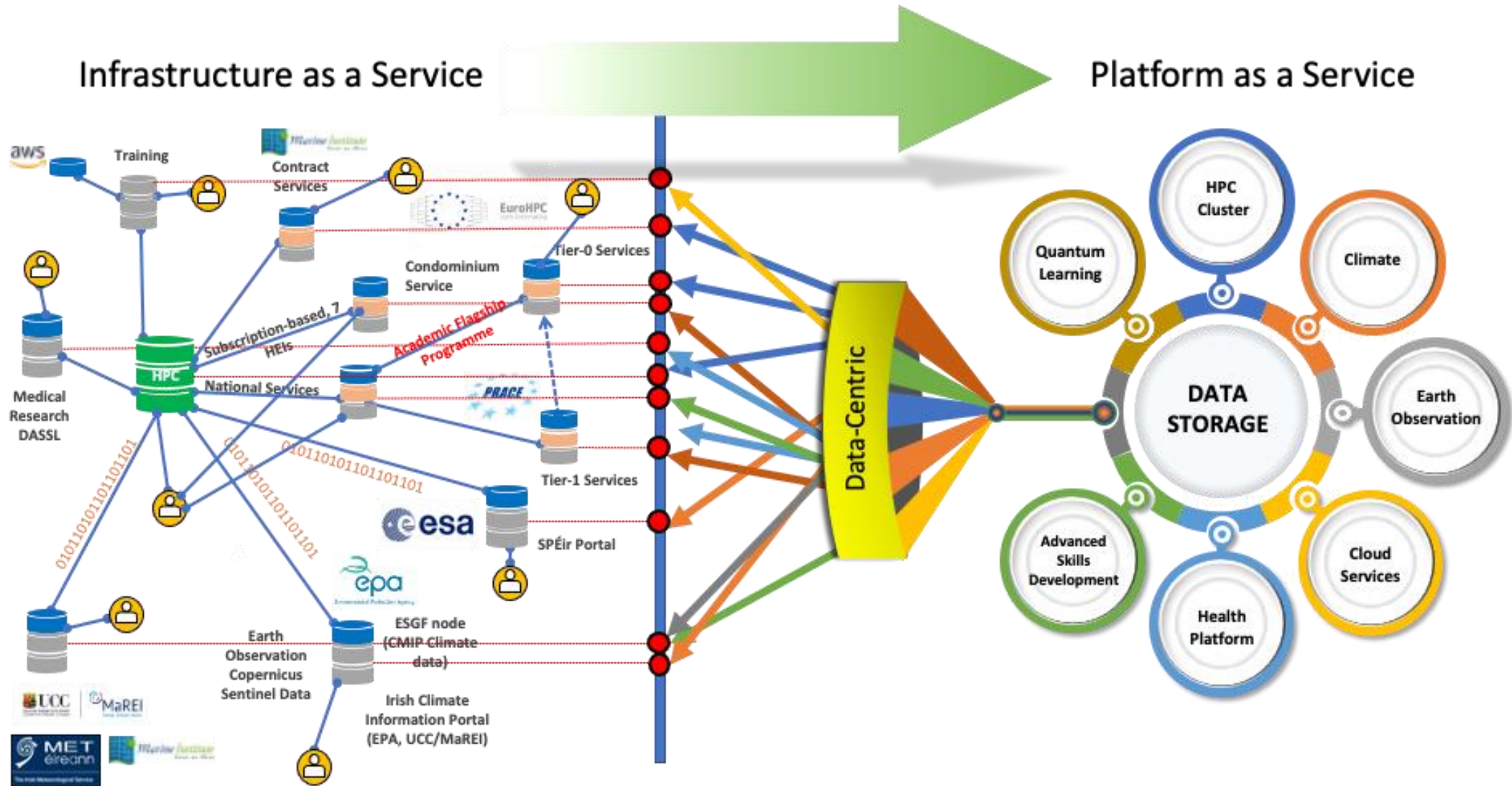
OLLSCOIL NA  
GAILLIMHÉ  
UNIVERSITY  
OF GALWAY



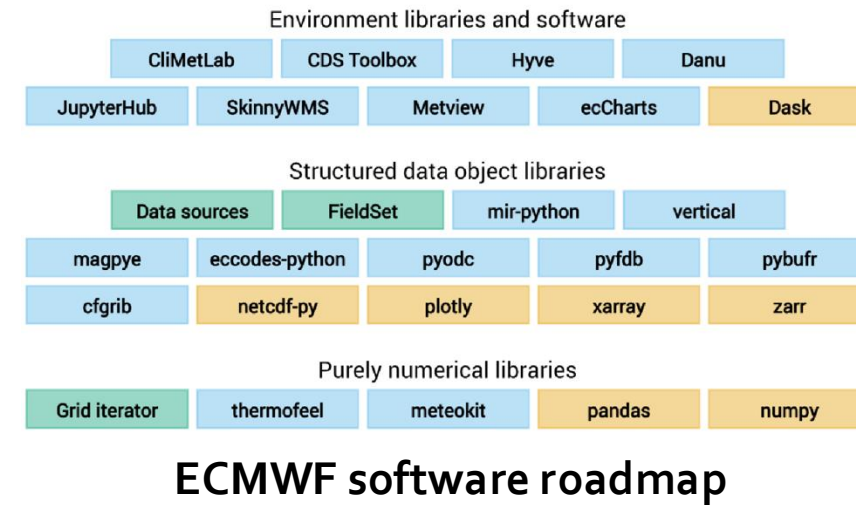
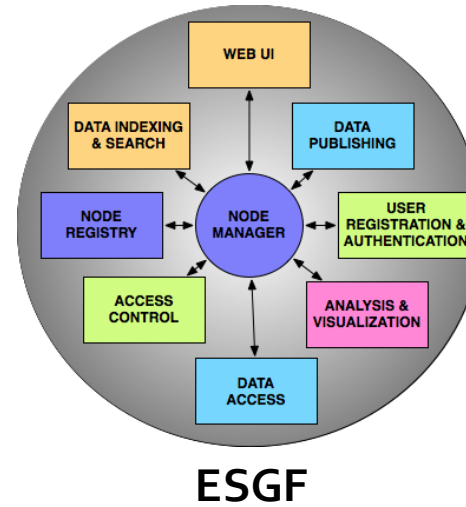
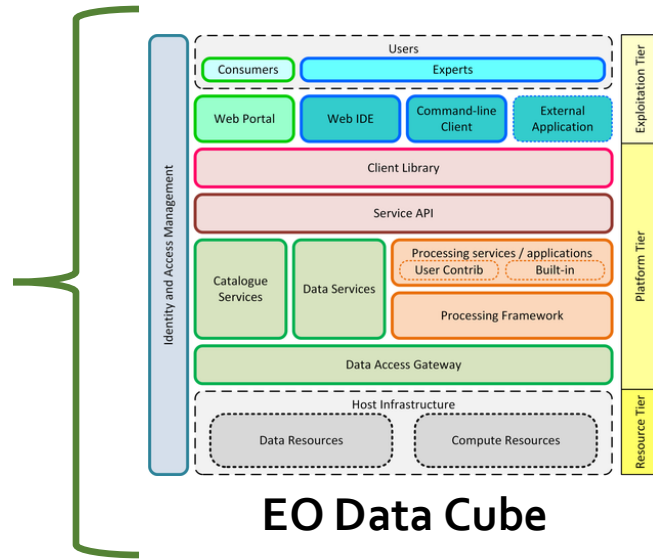
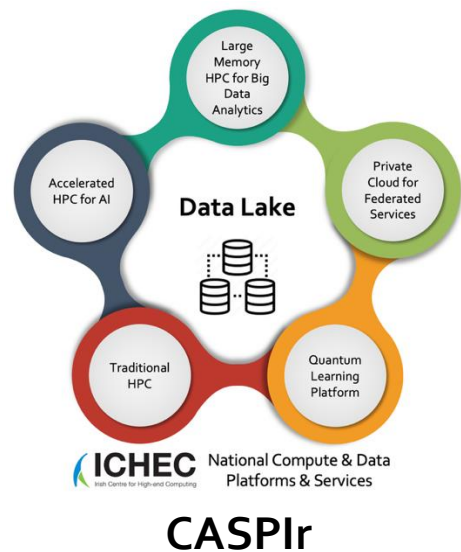
# Application / workflow driven Infrastructure Design

Venkatesh Kannan  
Irish Centre for High-End Computing (ICHEC)

# ICHEC Data & Compute Services



# ICHEC Data & Compute Services



## Operational Platforms & Services

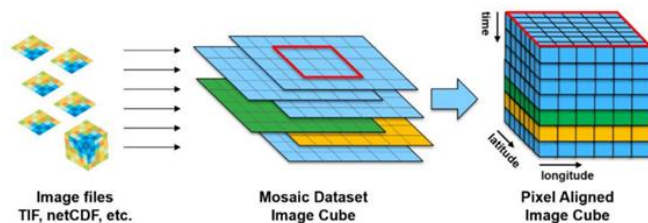


## Towards Data Spaces & Digital Twin Platforms



Special Area of Conservation (SAC)  
Protect & enhance High-Status Waters  
Results Based Payment Scheme (RBPS)

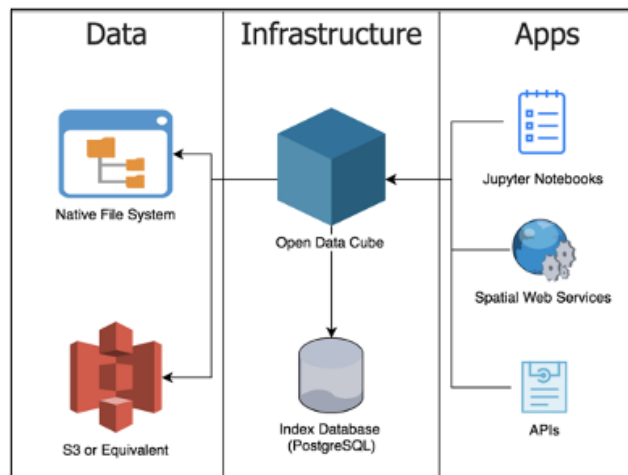
# Earth Observation Data Cube



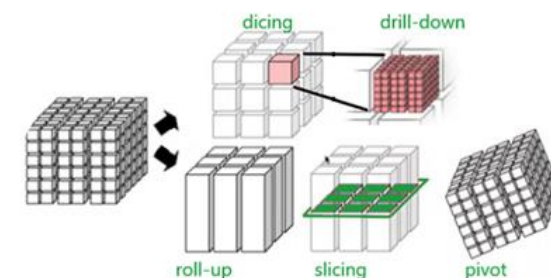
Prepare files for data cube



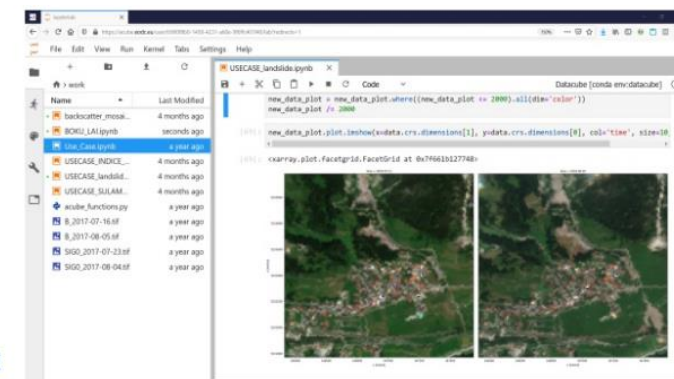
Run workflows  
continuously adding  
updates to data cube



Data served for  
HPC-enabled  
analytics / ML



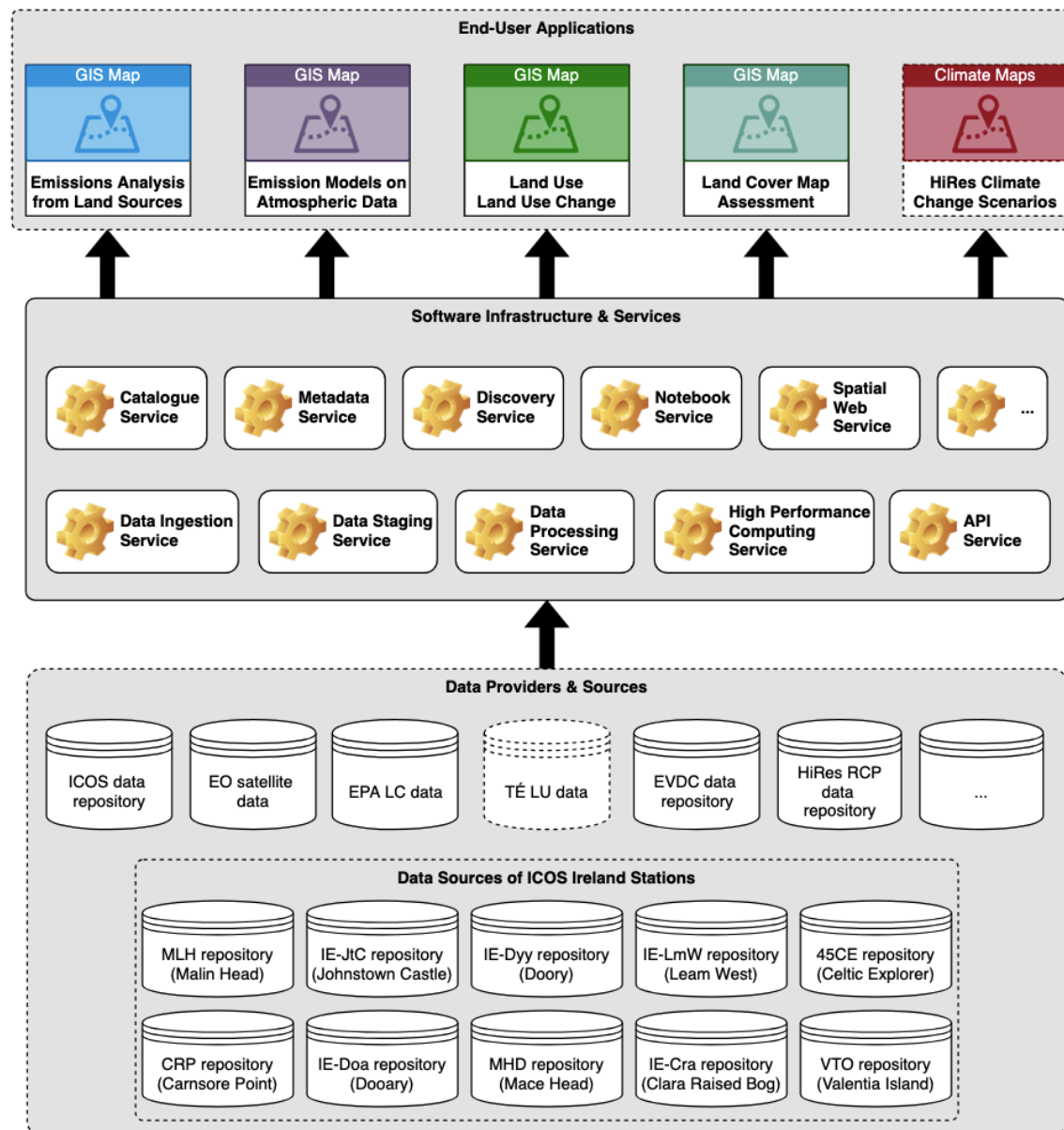
Map widgets in web pages,  
notebooks , ArcGIS, etc



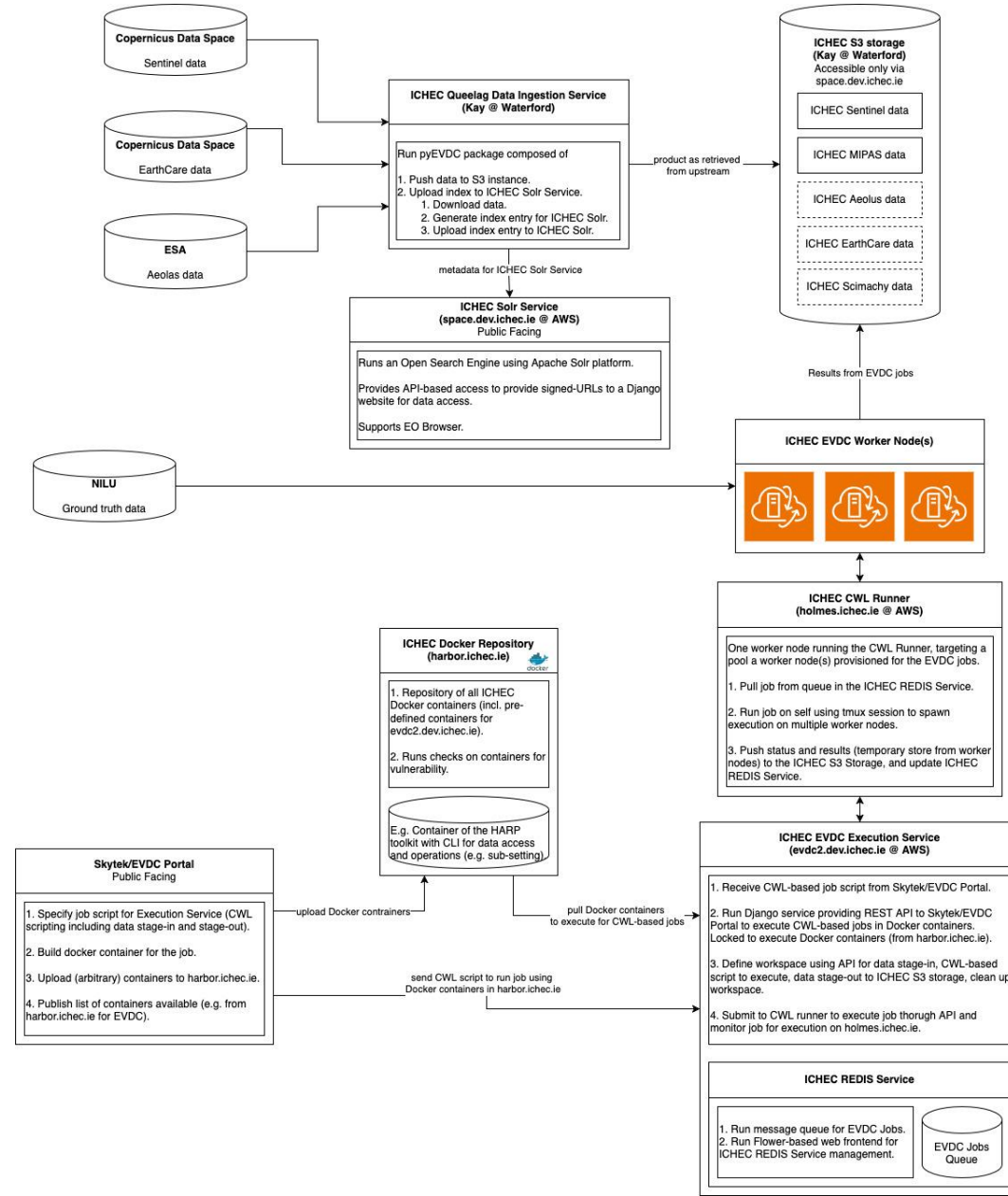




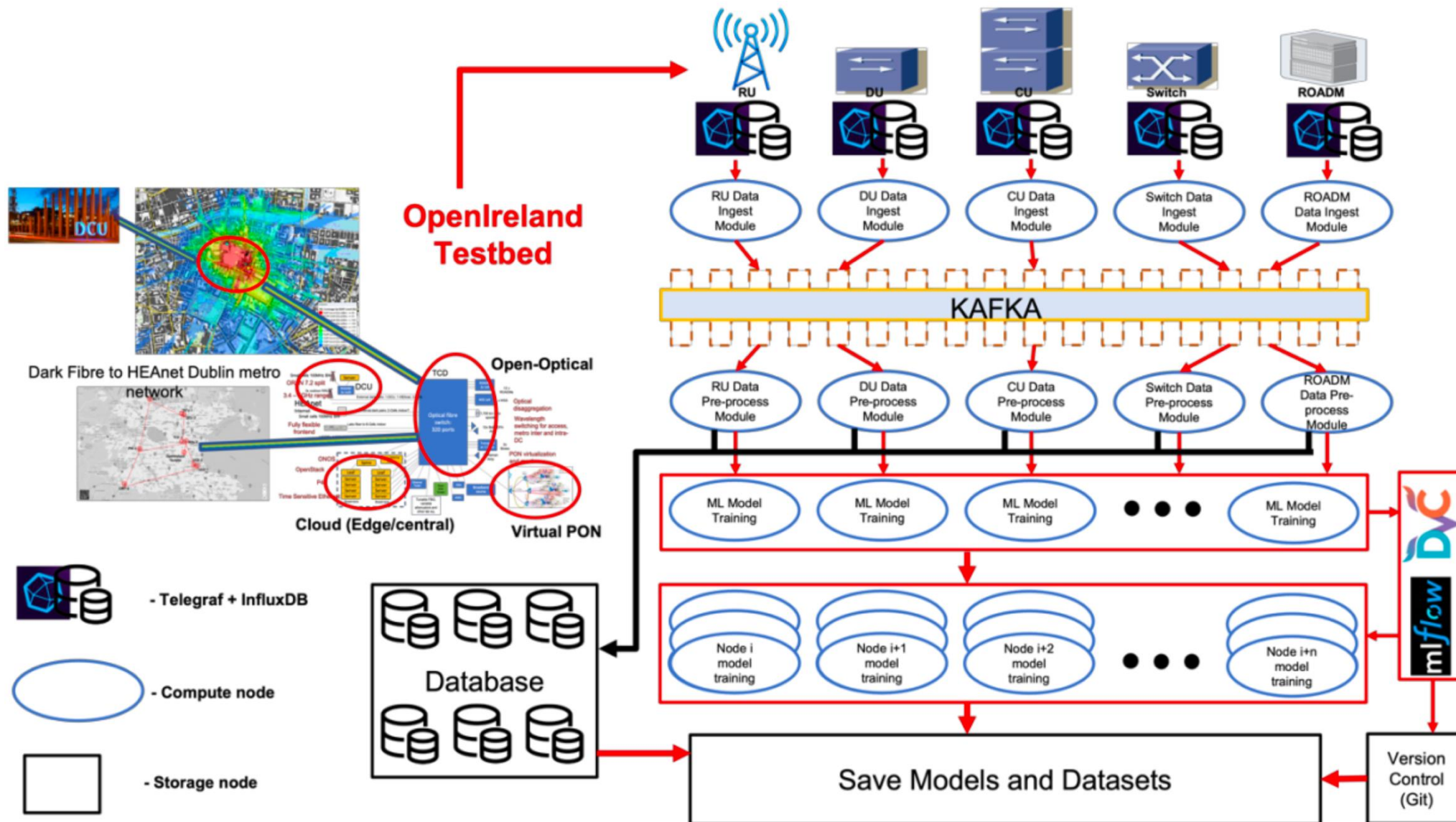
# GeoCube



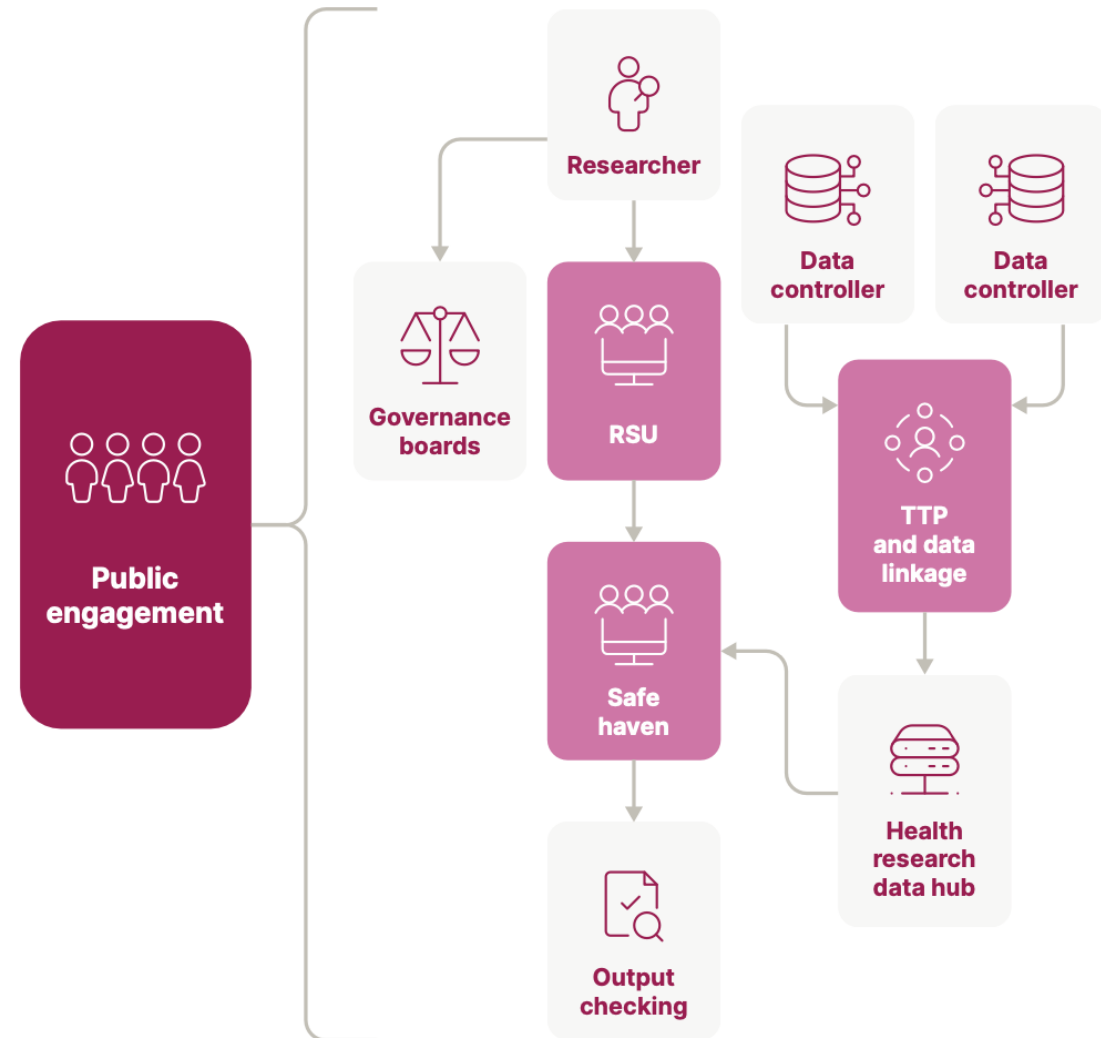
# ESA Validation Data Centre



# Telecommunication Testbed



# Trusted Research / Secure Data Environments



- Technical infrastructure for data management and processing
  - Secure data environments (SDE/TRE/VDR)
  - Health research data hub
  - ISO 27001 certification
  - Co-location with other research infrastructure vs. secure cabinets/cages
  - Others?
- Institutional governance (of ICHEC) to implement and operate technical infrastructure for health R&I data platforms, services
  - (E.g.) RSU/TTP/DLU/HDAB ↔ technical infrastructure provider
  - Handle DPIA, data input/management, data sharing agreements, repatriation of genomic data
  - Regulations/governance frameworks specific to themes (e.g., clinical research, cancer studies, genomics, etc.)
  - Others?



# Mapping data, compute, usage characteristics

## Data characteristics

- Average size (order of magnitude)
- Growth rate (frequency and size of accumulation)
- Ingestion rate (speed of new data ingested into system; recurrent, bulk?; throughput-critical?)
- Processing speed (based on use-cases on data; e.g., streaming, batch?)
- Data types and formats (e.g., satellite data, maps, images, tables, structured/unstructured, etc.)
- Read/write intensity (e.g., read-heavy, write-heavy, balanced?)
- Persistence (long-term preservation, warm/cold storage, purge after session, etc.)
- Data security (encryption, sensitivity, access control, authentication/authorization, etc.)
- Data flow patterns (e.g., ingest-curate-store, ingest-curate-analyse-store-results-purge-inputs, ingest-curate-transform-store-analytics, etc.)

## Compute characteristics

- Workload type (batch, real-time, interactive, hybrid)
- Intensity known (parallelism, memory requirement, acceleration, compute vs. IO)
- Input and output data characteristics (average per user job)
- Processing throughput (requests to handle per unit time, etc.)

## Usage characteristics

- Access mechanisms (CLI/GUI, batch/XaaS)
- Usage characteristics (EPA/Agencies/researchers, number of concurrent users, peak concurrent requests, burst traffic patterns)
- QoS and SLA (uptime, response time, certification, security compliance, data sensitivity, data redundancy/replication, etc.)
- Scaling estimations (horizontal vs. vertical, auto-scaling, etc.)

# Infrastructure Design Considerations

- Hardware
  - Compute | Data storage | Interconnects | Modularity
- Data Centre
  - Multi-site | Multi-tenancy | Tier-4
- Caging / Sandboxing
  - Software control | Data access control | Network restrictions | Execution isolation | Reproducibility
- Resource management / Scheduling
  - Workflow engines (DAG, job dependencies) | Heterogenous jobs | Hierarchical scheduling
- Software platform engineering
  - Containerisation | Heterogenous jobs | Security | Modularity | Reusability | Micro-services
- User access & management
  - Batch, interactive, automated workflows | Portals, notebooks, CLI | Quotas, on-demand, shared
- Monitoring
  - Hardware resource utilisation | Application workflows | Software platform components
- Governance
  - Data | Institutional

# IRL-DataSpaces

|                            |                       |
|----------------------------|-----------------------|
| Framework & Roles          | Architecture & Design |
| Security & Privacy         | Lifecycle management  |
| Compliance & Regulation    | Stewardship           |
| Ownership & Accountability | Metadata Management   |
| Access & Usage Policies    | Quality Management    |

Data Governance



|                          |                              |
|--------------------------|------------------------------|
| Security & Privacy       | Standards & Interoperability |
| Regulations & Compliance | Ecosystem & Services         |
| Lifecycle & Provenance   | Gravity & Performance        |
| Residency Requirements   | Open/Closed/Sensitivity      |
| Risk Management          | Producers, Users & Use-cases |

Sectoral Specifics



Technical Infrastructure

|                       |                               |
|-----------------------|-------------------------------|
| Resource Registration | Data Services                 |
| Catalogue & Discovery | VRE / Data Rooms              |
| Metadata Services     | HW & SW Infrastructure        |
| PID Services          | Certification & Cybersecurity |
| Data Ontology         | Federated AA(A)I              |

# IRL-DataSpaces

