WEBINAR

# Data spaces: Discovering the building blocks
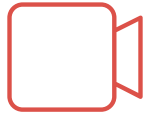
data.europa academy

# Rules of the game

The webinar will be recorded

For questions, please use the ClickMeeting chat.

Please reserve 3 min after the webinar to help us improve by filling in our feedback form

# Introduction

**Flora Kopelou**
**Data.europa.eu,**
**Publications Office of the**
**EU**

**Clara Pezuela**
**FIWARE Foundation**

**Edward Curry**
**University of Galway**

**Alexandra Balahur**
**European Commission**
**DIGIT**

# Agenda

| | |
|---|---|
| 10.00 – 10.10 | Opening and introduction of the series |
| 10.10 – 10.30 | The building blocks of data spaces – *Clara Pezuela* |
| 10.30 – 10.50 | Exploring the Research  Challenges with Data Spaces – *Edward Curry* |
| 10.50 – 11.10 | SEMIC Specifications, services and trainings in support of Data Spaces – *Alexandra Balahur* |
| 11.10 – 11.25 | Questions and answers |
| 11.25 – 11.30 | Closing remarks |

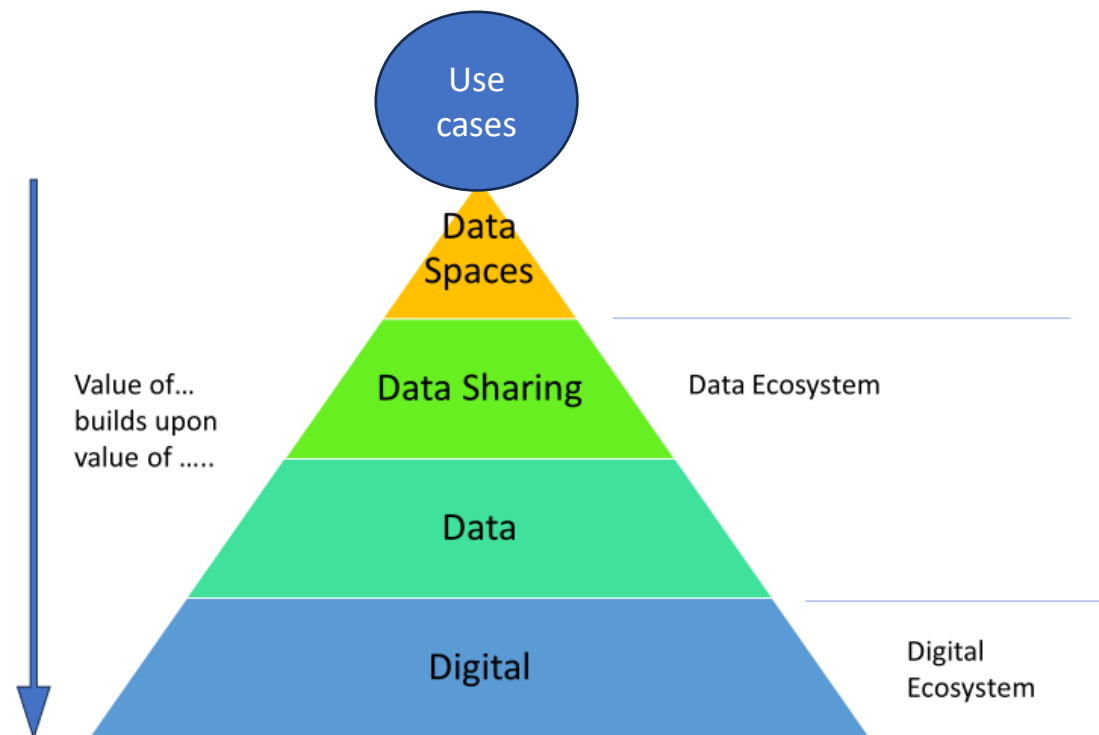# The building blocks of data spaces

Clara Pezuela

# Discovering the building blocks

Clara Pezuela (FIWARE Foundation)

Data Europe webinar - 6 October 2023

# Context of data spaces

- The value of data spaces relies on the value of data sharing and therefore on the value of data

- Data Spaces operate in digital ecosystems and inherit some of their characteristics

- Data Spaces were introduced in the European Data Strategy to support the idea of "single market of data"

- Being concrete through use cases



Use cases

Data Spaces

Data Sharing

Data

Digital

Value of... builds upon value of .....

Data Ecosystem

Digital Ecosystem

# Data Space concept

> *A distributed system defined by a governance framework, that enables trustworthy data transactions between participants while supporting trust and data sovereignty (DSSC Glossary)*

## What is?

- Several infrastructures supporting one or more use cases:
    - Distributed structure
    - Governance framework
    - Enabling trusted data transactions
    - Enabling data sovereignty

## What is not?

- A data lake
- Only a data platform
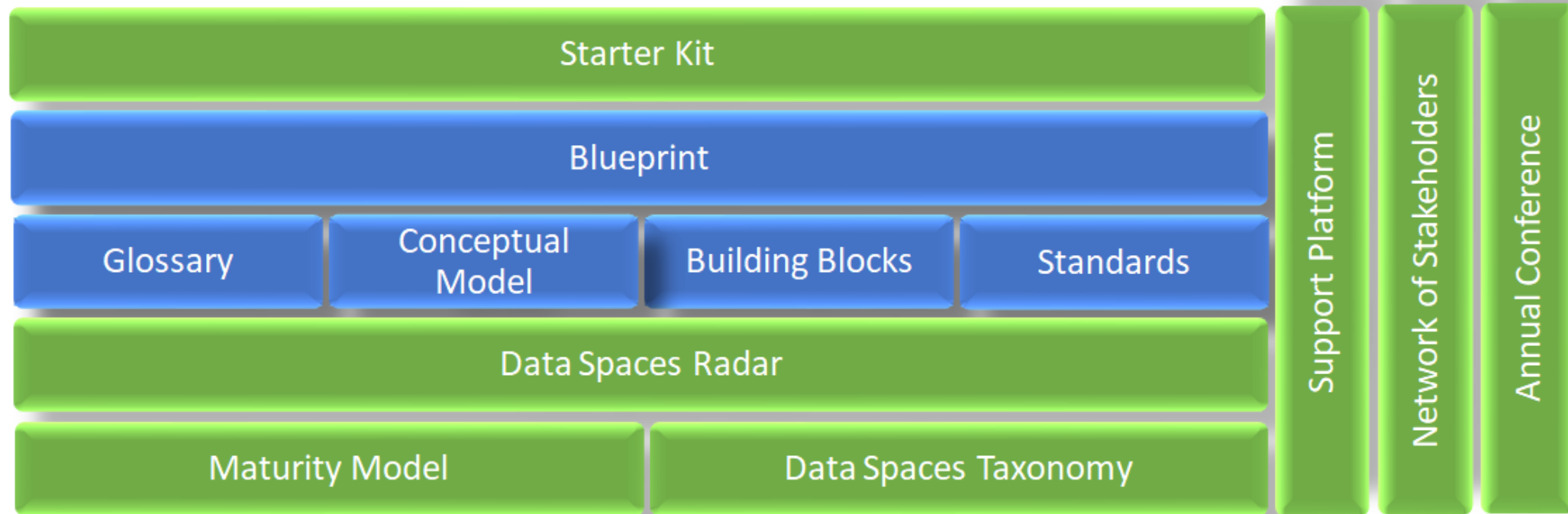- Only a digital ecosystem

# DSSC mission

- Enable data spaces to reach a higher flight level faster: a **quick start** and an **accelerated scale-up**

- Provide the tools to address the **basic organizational and technical matters**, required by every data space

- This includes a **blueprint**, best practices, common standards and reference implementations which will be developed according to a **co-creation process**

- Enable dataspaces to focus on their domain-specific business challenges and **provide business benefits to their participants**

# DSSC assets



DELIVERY PLAN

# Blueprint v0.5 just released

https://dssc.eu/page/knowledge-base



Starter Kit    Blueprint    Building Blocks    Conceptual model    Glossary

# Blueprint

# Data Spaces Blueprint Structure

# Blueprint v0.5 elements

## Glossary

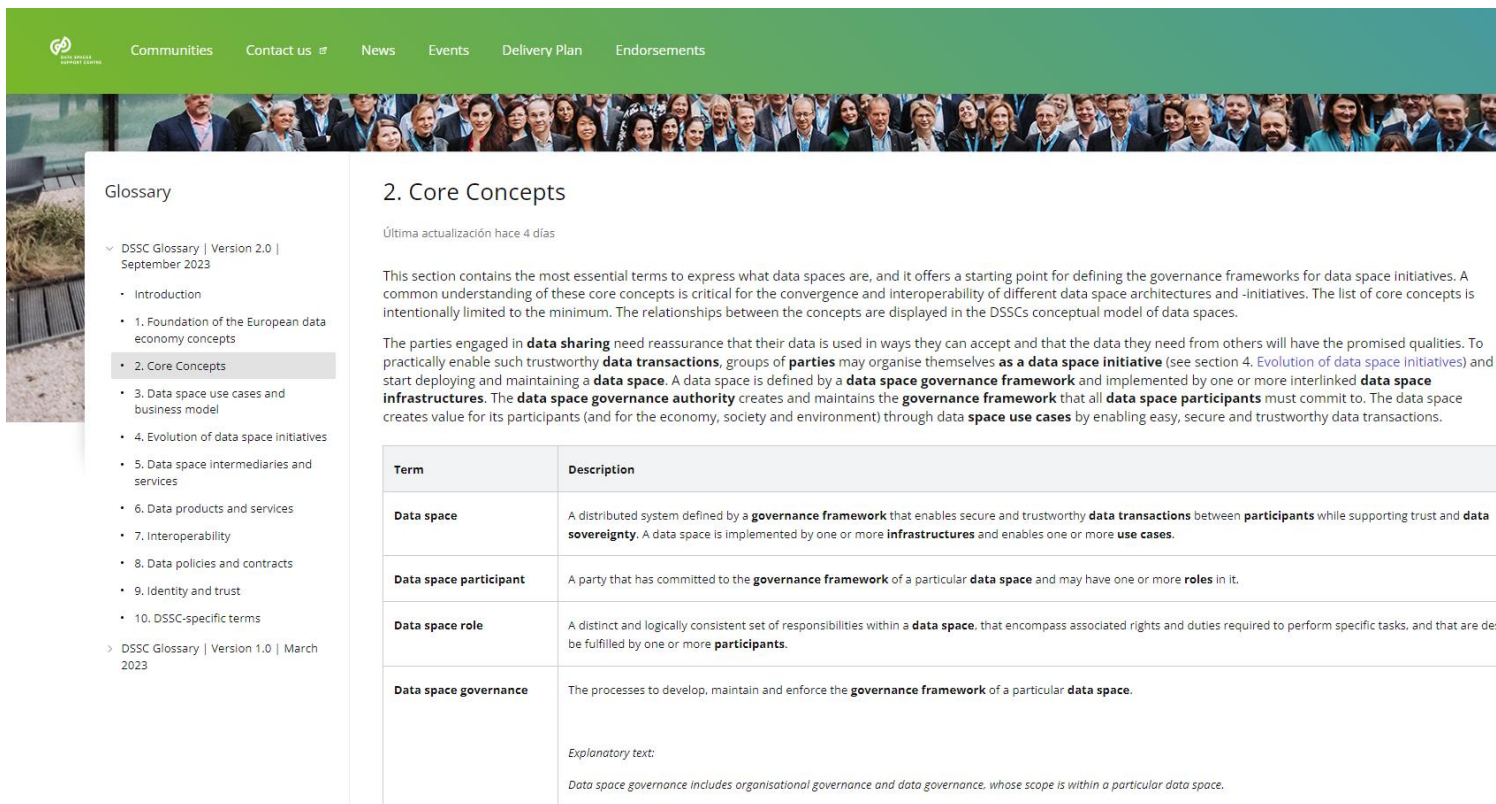- Curated set of terms and definitions
- V2.0

## Conceptual Model

- Model which represents concepts and relationships among them
- V0.5

## Building Blocks

- Taxonomy of basic units or components that can be implemented and combined with other building blocks to achieve the functionality of a data space
- V0.5

# Glossary

- In co-development with data spaces

- Common vocabulary for DSSC publications and communications

- Terms are provided with a criterion and definition
  - Same term may have different meanings in different contexts

- Others may keep their own terminology

- Adopting definitions from EU legislation and international standards, when possible

# Conceptual Model

- Provide a well-defined language
- Provide a high-level view of what a data space is
- Future structure will provide:
  - An overarching view of the data space environment (level 0)
  - Detailed view on individual concepts (level 2)

Level 1: basic terminology and key concepts

# Building Blocks (BBs)

**Organisational & business building blocks**

Describe capabilities on a governance, legal and business level.

**Technical Building Blocks**

Describe capabilities on a technical level about:

*Data Interoperability*
*Data Sovereignty & Trust*
*Data Value Creation*

## Organisational and Business Building Blocks

| Business | Governance | Legal |
|---|---|---|
| Business Model Development | Organisational Governance | Regulatory Compliance |
| Use Case Development | Data Sharing Governance | Contractual Framework |
| Data Product Development | | |
| Data Space Intermediary | | |

| Data Interoperability | Data Sovereignty & Trust | Data Value Creation |
|---|---|---|
| Data Models | Access & usage policies and control | Data, Services and Offerings descriptions |
| Data Exchange | Identity Management | Publication and Discovery |
| Provenance & traceability | Trust | Marketplaces |

## Technical Building Blocks

# Organisational and Business BBs - Business

- Provides the essential concepts necessary in the business modelling of a data space

- Data space business model vs data space participant business model

| | |
|---|---|
| **Business model development** | • Define the data space business model and identify some considerations the governance authority should take into account while developing its business model |
| **Use case development** | • Strategic approach to amplify the value of a data space by fostering the creation, support and scaling of use cases |
| **Data producto development** | • Data product templates for the data providers, data space governance rules for the data products and network effects between data providers and users. |
| **Data space intermediaries** | • How to make business and governance decisions related to data space intermediaries |

# Organisational and Business BBs - Governance

- Focused on data space-level governance, emphasizing its dynamic nature.

- Governance in data spaces needs to adapt as the data space evolves.

- This requires data space participants in the data space to work together strategically for effective governance.

**Organisational governance**
- Guides setting up the data space governance authority by identifying key decision points and options for establishing inclusive governance and transparent rules and roles.

**Data sharing governance**
- Supports the governance authority in establishing common rules that promote effective and reliable data sharing processes and introduces different ways to organise data transactions within a data space.

# Organisational and Business BBs - Legal

DATA SPACES
SUPPORT CENTRE

- Provide guidance and resources for the data space initiatives to ensure compliance with legislation and establish a robust contractual framework.

**Regulatory compliance**

- Awareness of the legal landscape and assessing applicable regulatory requirements to ensure legal compliance and alignment with EU values

**Contractual framework**

- Establishing clear and enforceable rights and obligations for data space participants and provides contractual resources for data space participants to regulate their data transactions.

Funded by
the European Union

# Technical BBs – Data Interoperability

| | Description | Key capabilities | Link to specifications |
|---|---|---|---|
| Data Models | Capabilities to define and use shared semantics in a data space | • Vocabularies<br>• Vocabulary management process<br>• Vocabulary hub | Data Models |
| Data Exchange | Capabilities relating to the actual exchange and sharing of data | • Meta specifications and best practices for the adoption of existing data exchange APIs<br>• Generic purpose data exchange API, including methods to query, update and delete data<br>• Tooling to use and maintain data exchange APIs | Data Exchange |
| Provenance & Traceability | Capabilities for tracking the process of data sharing, so it becomes traceable and compliant | • Framework for requirements for observability<br>• Third parties to provision or use evidence<br>• Mechanisms to provide and use evidence of the activities of a transaction | Provenance & Traceability |

Funded by
the European Union

# Technical BBs – Data Sovereignty & Trust

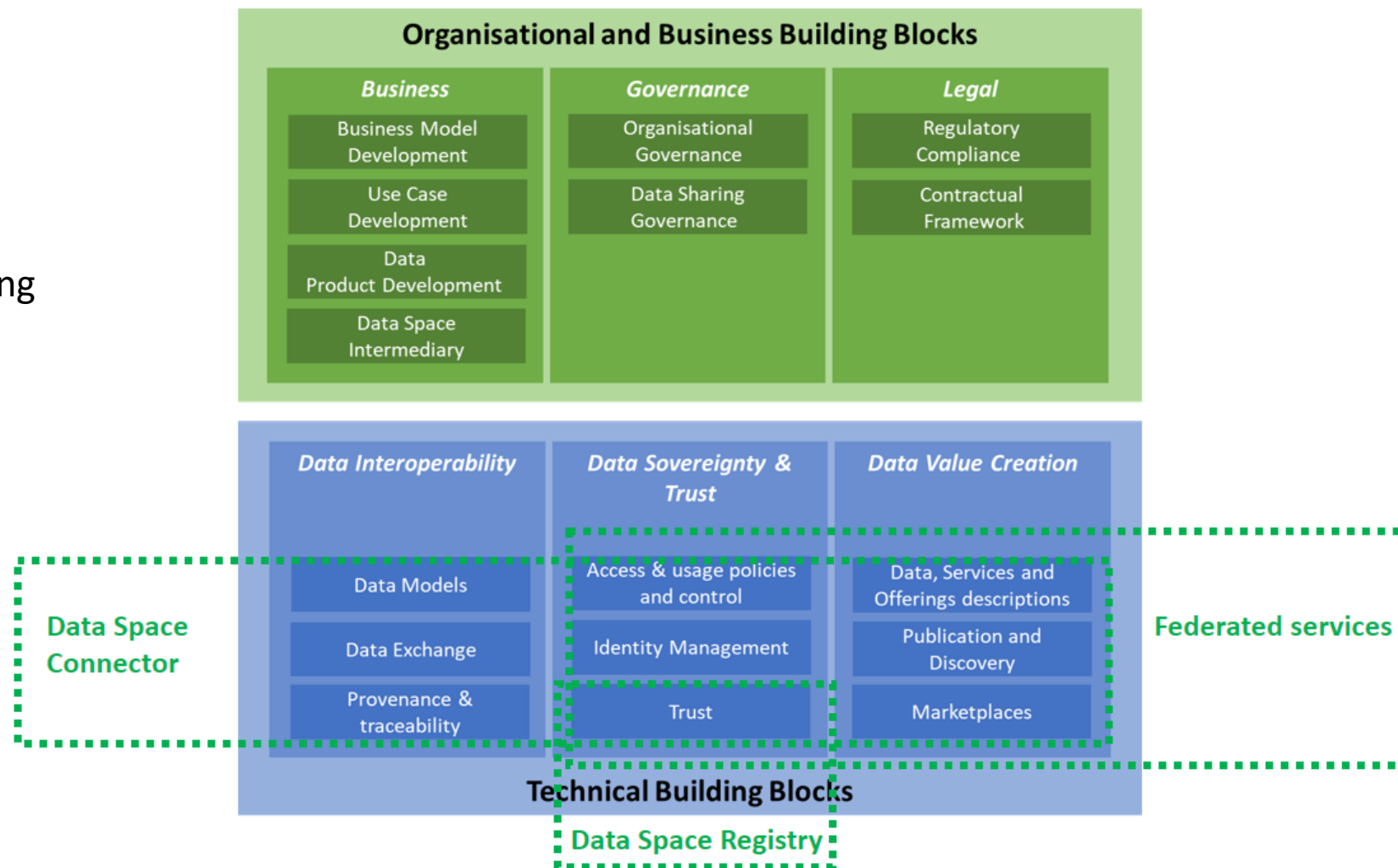| | Description | Key capabilities | Link to specifications |
|---|---|---|---|
| Access and usage policies and control | Ability to specify and policies within a given data space, by the data space authority and the individual participants | • Access policies specify the conditions to access services and data<br>• Usage Policies specify rights and obligations for the usage of the data, this includes future usages of data<br>• To enable the decision-making process in evaluation policies connection to other building blocks is required for Identification, Claim Management, Authentication and Authorization | Access and usage policies and control |
| Identity Management | Management of identities within a data space | • Onboarding (and offboarding) participants into a data space, by providing them with an identity of verifying their identity<br>• Issuing, holding and verifying identities | Identity Management |
| Trust | Being able to verify that a participant of a data space adheres to certain rules | • Semantic to describe and verify the description of service, data products and policies, notably for access control, usage purpose, consent, authorisation and rights delegation<br>• Level of assessment activity - either a declaration or certification - for making specific claim<br>• List of accredited parties eligible to issue identifiers, certifications and cryptographic material such as certificates<br>• Procedure for parties to request and issue identifiers for themselves, services or data products | Trust |

# Technical BBs – Data Value Creation

| | Description | Key capabilities | Link to specifications |
|---|---|---|---|
| Data, Services and Offerings descriptions | provides to data providers the tools to describe appropriately, and in a complete way, a data product, in a manner that will be understandable by any participant in the data space. | • detailed explanation assists users in understanding the data product content and structure<br>• rich description allows users to get an impression of the data product before conducting any analysis or processing<br>• complete data description is required for evaluating the dataset's quality and dependability<br>• well-documented data description contributes to data governance activities such as data privacy, security, and regulatory compliance | Data, Services and Offerings descriptions |
| Publication and Discovery | Allows data providers to publish the description of their data, services and offerings, in a way that they can be discovered by users, following the FAIR principles | • Management of self-descriptions<br>• Discovery of self-descriptions<br>• Enable dynamic transactions<br>• Mange access to self-descriptions | Publication and Discovery |
| Marketplaces | Marketplace capabilities, in such a way that provider and user can enter into a relationship about the access, provision and use of a data product | • Catalogue management<br>• Contract and producto order management<br>• Aftersales assistance and customer service | Marketplaces |

# Systems architecture view

- Technology is needed to implement the capabilities of Technical BBs
- Certain services are implementing the BBs
- To facilitate the onboarding of data space participants



**Organisational and Business Building Blocks**

| Business | Governance | Legal |
| --- | --- | --- |
| Business Model Development | Organisational Governance | Regulatory Compliance |
| Use Case Development | Data Sharing Governance | Contractual Framework |
| Data Product Development | | |
| Data Space Intermediary | | |

**Technical Building Blocks**

| Data Interoperability | Data Sovereignty & Trust | Data Value Creation |
| --- | --- | --- |
| Data Models | Access & usage policies and control | Data, Services and Offerings descriptions |
| Data Exchange | Identity Management | Publication and Discovery |
| Provenance & traceability | Trust | Marketplaces |

Data Space Connector

Federated services

Data Space Registry

# How to use the Blueprint

- When posible follow the proposed structure in your data space

- Do the mapping between your building blocks and DSSC ones

- Identify missing or incomplete functionalities

  - Specific for your domain – refer to the DSSC building block that you are extending in your domain

  - Relevant for other domains – propose the addition to the DSSC to be incorporated to the DSSC taxonomy

- Keep updated about future releases

- Propose implementations (upcoming)

# Summary

- Demand on data sharing to address business and societal challenges

- Data Spaces as instruments to facilitate the data sharing

- Bringing an extra value to the existing ways of data sharing

- Guidelines and tolos are needed for data spaces initiatives

- DSSC is here to support the understanding and development of data spaces

- Providing a set of assets, including a blueprint

- Blueprint v0.5 is out! Go and take it!

# How to keep posted

- News and publications at dssc.eu
  - Get support

- Join the Thematic Groups
  - Knowledge base and periodic meetings
  - Get involved

- Monthly DSSC Insights Series
  - Next one about the Blueprint on 12th October at 16:00

- Data Spaces Symposium in March 2024

Funded by
the European Union

# Thanks!

More info at contact@dssc.eu

**Funded by
the European Union**

**DATA SPACES
SUPPORT CENTRE**

# Exploring the Research Challenges with Data Spaces

Edward Curry

Common European Data Spaces

# Foundations of the Web go back to the 60s….

The first ever internet message is sent
**1969**

The Domain Name System is born
**1983**

World.std.com becomes the first commercial provider of dial-up access
**1989**

The World Wide Web goes mainstream
**1991**

**1965**
The first ever WAN (Wide Area Network) is established

**1970-1976**
The rise of the LAN (Local Area Network)

**1987**
Cisco ships its first router

**1990**
Tim Berners-Lee invents HTML

**2018**
4 billion internet users around the globe

Insight

# Common European Data Spaces ……..A Ten Year Journey

# Research Challenges

E. Curry, S. Scerri, and T. Tuikka, "Data Spaces: Design, Deployment, and Future Directions," in *Data Spaces*, 2022.
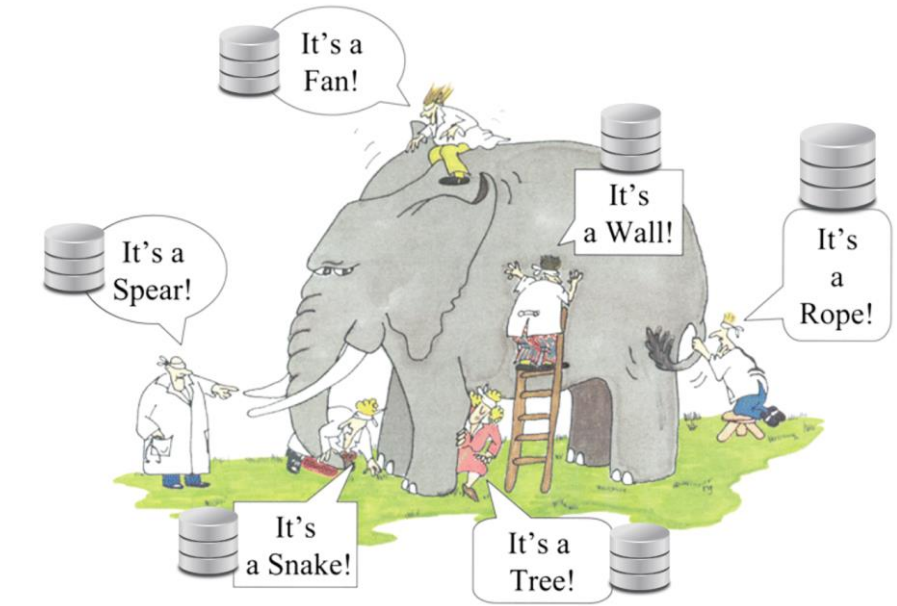
# Expanding Content Space...

- Heterogeneous, complex and large-scale data

- Very-large and dynamic "schemas"

- Open Environments: distributed, decentralised decoupled data sources, anonymous users, multi-domain, lack of global order of information flow

circa 2000

10s-100s attributes





circa 2020

1,000s-1,000,000s attributes

- Multiple perspectives (conceptualisations) of the reality.

- Ambiguity, vagueness, inconsistency.



... Fundamental Decentralisation

Insight

# Increasing Amounts of Multimodal Data....



(More than 220k people in Ireland are blind or visually impaired)

(An average of 110 drownings occur in Ireland every year)

(Every 3 minutes in Ireland someone gets a cancer diagnosis; every hour someone dies from cancer)

.J. Khan, J.G. Breslin, E. Curry, "Expressive Scene Graph Generation using Commonsense Knowledge Infusion for Visual Understanding and Reasoning," Extended Semantic Web Conference (ESWC) 2022, Crete, Greece, May-June 2022.

Insight

# Increasing amounts of Subjective Data...

"The very concept of objective truth is fading out of the world."

*George Orwell*

Insight

# Subjective and Objective Attributes and Query



Table 3: Subjective attributes in different domains.

| Domain | %Subj. Attr | Some examples |
|---|---|---|
| Hotel | 69.0% | cleanliness, food, comfortable |
| Restaurant | 64.3% | food, ambiance, variety, service |
| Vacation | 82.6% | weather, safety, culture, nightlife |
| College | 77.4% | dorm quality, faculty, diversity |
| Home | 68.8% | space, good schools, quiet, safe |
| Career | 65.8% | work-life balance, colleagues, culture |
| Car | 56.0% | comfortable, safety, reliability |

 Li, Y., Feng, A., Li, J. *et al.* Querying subjective data. *The VLDB Journal* **30,** 115–140 (2021). https://doi.org/10.1007/s00778-020-00634-5

Insight

# The Red Queen Hypothesis

*"It takes all the running you can do, to keep in the same place. If you want to get somewhere else, you must run at least twice as fast as that!"*

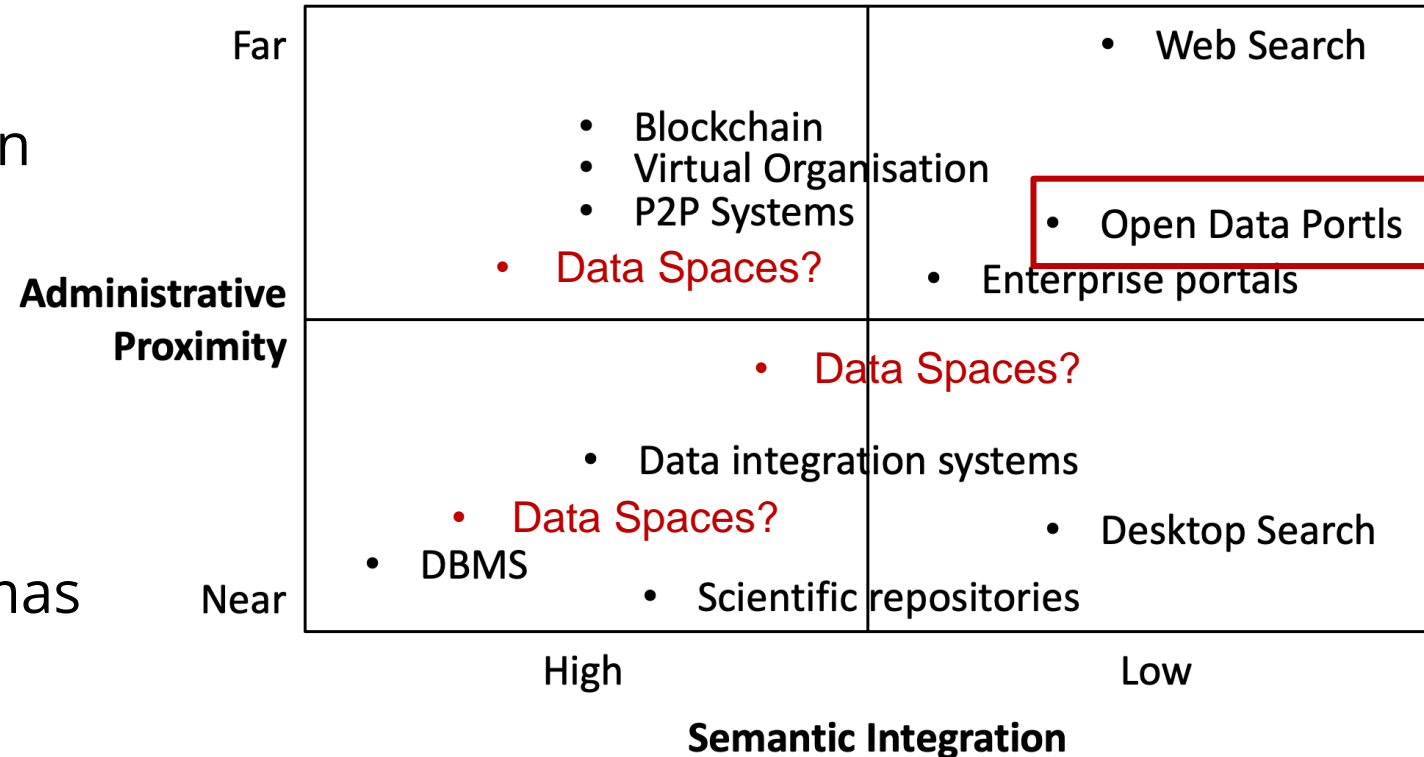Lewis Carroll's Through the Looking-Glass

# Control and Coordination

**Administrative Proximity**
- Close vs. Loose Coordination
- Assumptions concerning guarantees such as data, access, quality, and consistency,
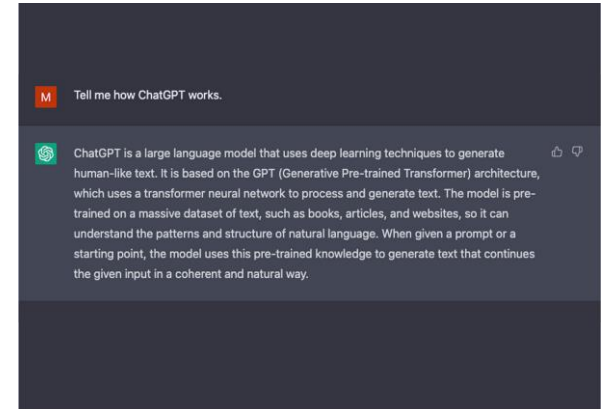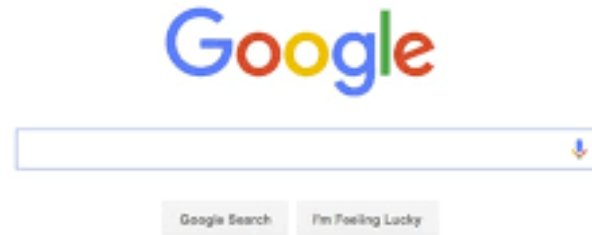
**Semantic Integration**
- Degree to which data schemas are matched up (types, attributes, and names).



Halevy, A., Franklin, M. and Maier, D. 2006. Principles of dataspace systems. *25th ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems - PODS '06* (New York, New York, USA, 2006), 1–9.

# Evolving Human Interactivity...



**From Structure** → **to Search** → **to Knowledge Graph** → **to Conversations ?**

~1995
~100K Websites
Exact Results
Human Curated

~1998
~2.4M Websites
Approximate Results
Computed

~2012
~700M
Approximate Results + Exact
Computed + Crowd

~2023
~2B
Approximate Results ?
?

Article development led by **acmqueue**
queue.acm.org
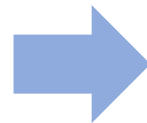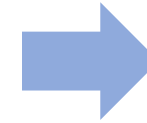
**In today's humongous database systems, clarity may be relaxed, but business needs can still be met.**

BY PAT HELLAND

# If You Have Too Much Data, then 'Good Enough' Is Good Enough

*"We can no longer pretend to live in a clean world. SQL and its Data Definition Language (DDL) assume a crisp and clear definition of the data, but that is a subset of the business examples we see in the world around us. It's OK if we have lossy answers—that's frequently what business needs."*

Insight

# What is a Data Space?

"Dataspaces are not a data integration approach; rather, they are more of a **data co-existence approach**. The goal of dataspace support is to provide base functionality over all data sources, regardless of how integrated they are." (Halevy, A., Franklin, M. and Maier, D. 2006.)

Incrementalism, Approximate, Interactive

Insight

# Data Space Enablers

**Data Space Support Platform** (Halevy et al.)

- Must deal with **many different formats** of data.
- **Does not subsume** existing systems; they still provide individual access via their native interfaces.
- Queries in are provided on a **best-effort and approximate basis**.
- Must provide **pathways to improve the integration** among the data sources, including streams and events, in a pay-as-you-go fashion.
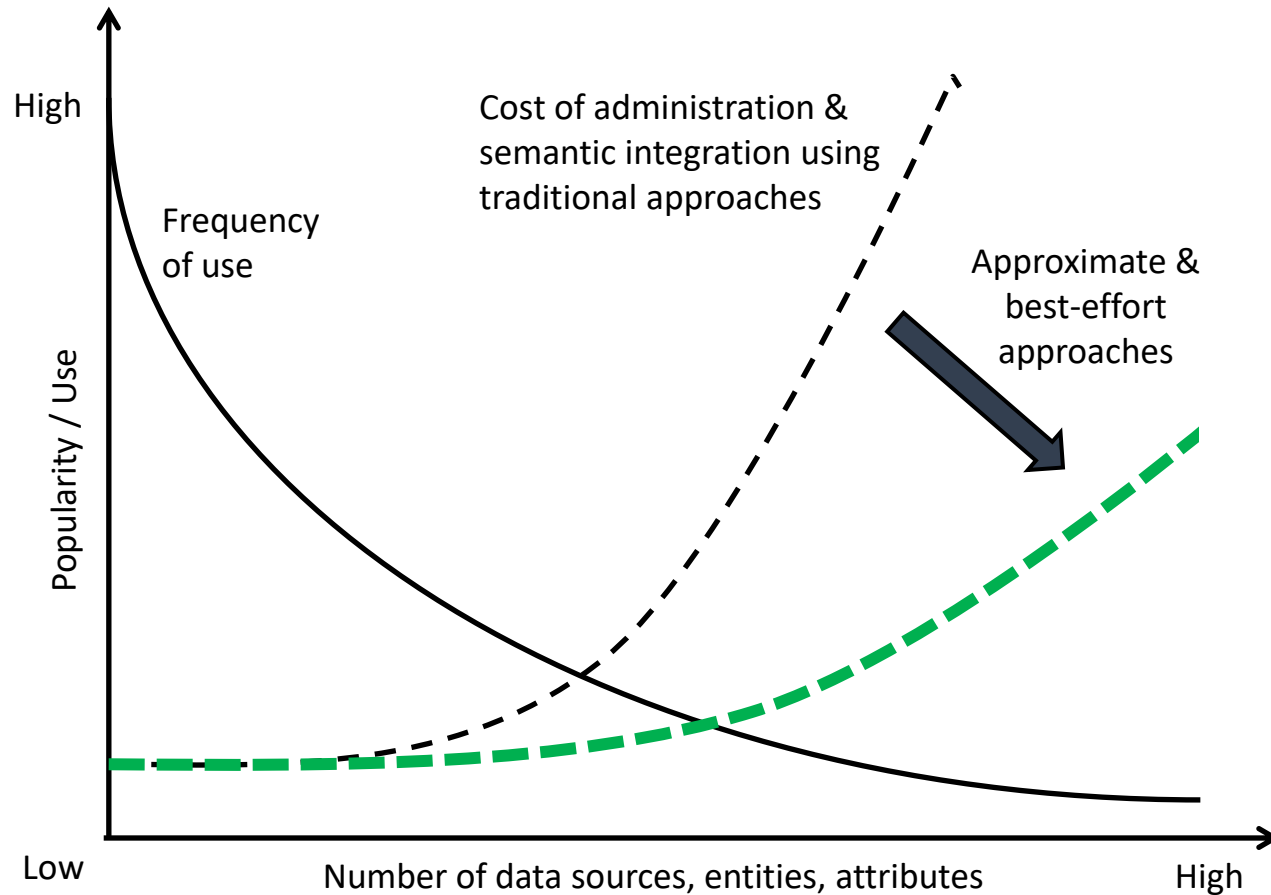
**Data Space Support Services**

- Catalog and Browse
- Search and Query
- Local Store and Index
- Discovery
- Enhancement
- Administration

## Reusing Human Attention

- Learn from users' activities
  - Create meaningful relationships between data sources
  - Enhance data sources

# Approximate and Best Effort Approaches
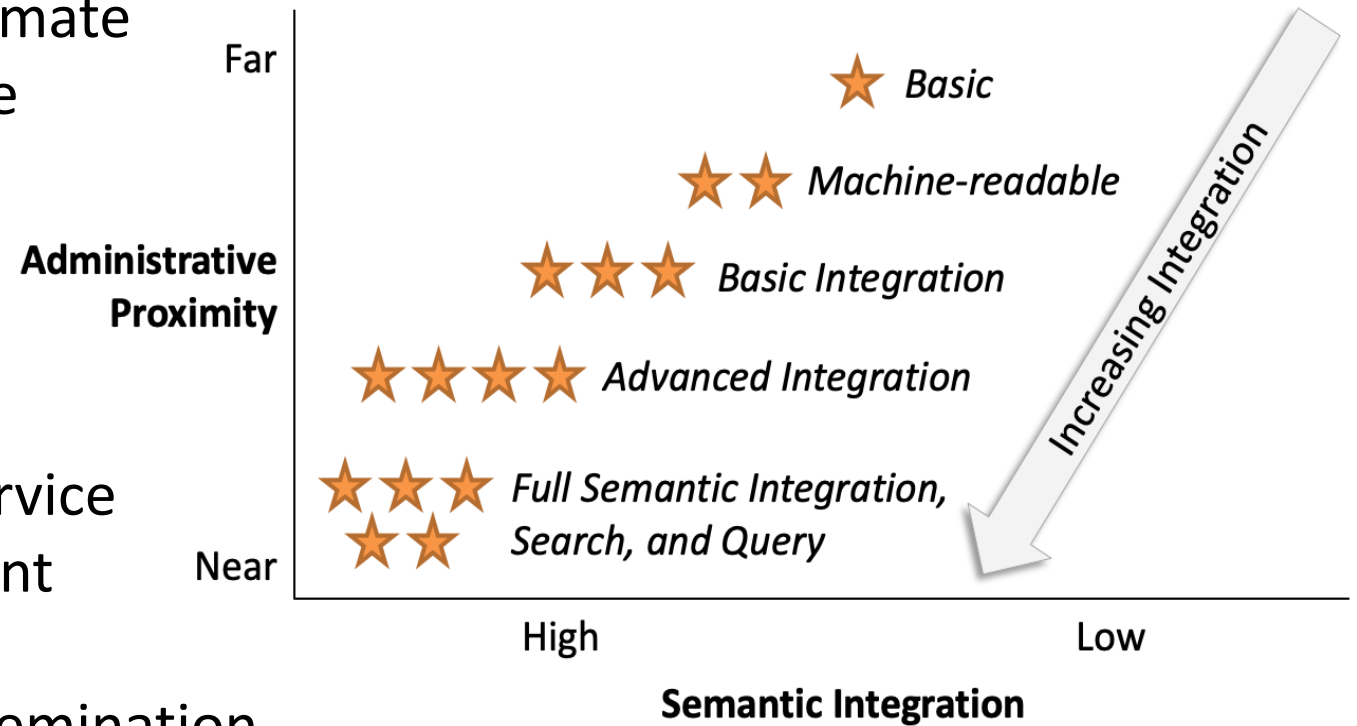


**The Long Tail of Data**

Insight

# Creating Approximate Services to support incremental data management is a key challenge…

Investigate techniques to enable approximate and best-effort support services for loose administrative proximity and semantic integration

**Incremental support services**

- Catalog
- entity management
- query and search
- data discovery
- human tasks

- quality of service
- complex event processing
- streams dissemination
- approximate semantic event matching

Far

**Administrative Proximity**

★ Basic

★★ Machine-readable

★★★ Basic Integration

★★★★ Advanced Integration

★★★
★★ Full Semantic Integration, Search, and Query

Near

High                    Low

**Semantic Integration**

*Increasing Integration*

**Real-time Linked Dataspaces**

OPEN ACCESS

# Research Trends

# Hybrid Neuro-Symbolic Approaches

- We need Semantics, not just statistics!

- Combing rules-based AI approaches (Knowledge Graphs) with statistical learning techniques (deep learning)

**Neuro-Symbolic AI**

**Neural Networks**
*(Statistical Learning + Reasoning)*

**Knowledge Graphs**
*(Expert Knowledge + Reasoning)*



**+**

Insight

# Examples of Symbolic Knowledge



Fig. 8. Illustrative overview of symbolic knowledge representations in NeSy.

- Neural- Symbolic Integration
- Knowledge Representation
- Knowledge Embedding

Insight

# Examples of NeSy Systems
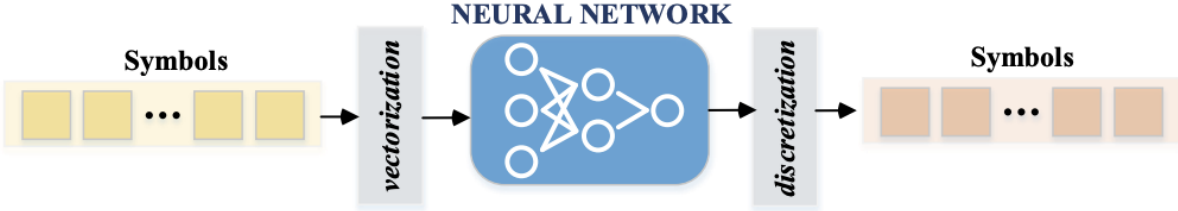


Fig. 2. Type 1: Symbolic Neuro Symbolic.
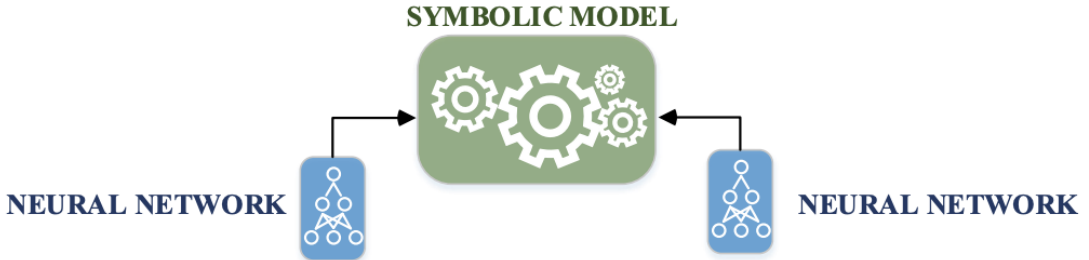


Fig. 3. Type 2: Symbolic[Neuro].



Fig. 4. Type 3: Neuro|Symbolic.



Fig. 5. Type 4: Neuro:Symbolic→Neuro.

Foundational Model

Wang, Wenguan, and Yi Yang. "Towards data-and knowledge-driven artificial intelligence: A survey on neuro-symbolic computing." *arXiv preprint arXiv:2210.15889* (2022)

Insight

# What is a Foundation Model?

Insight

# Definition

- *Any model that is trained on broad data (generally using self-supervision at scale) that can be adapted (e.g., fine-tuned) to a wide range of downstream tasks*

- Foundational models can be applied to a single modality while others are multimodal
  - BERT, GPT-3, GPT-4, LLaMA, Segment Anything Model, CLIP and Stable Diffusion

- **Emergence:** the behaviour of a system to be implicitly induced rather than explicitly constructed

- **Homogenization:** the consolidation of methodologies for building machine learning systems across a wide range of applications

Insight

# Centralize information from various modalities, then adapt to a wide range of downstream tasks....



- Trained on extensive datasets to establish a knowledge foundation

- Leverage transfer learning and scale in various downstream tasks.

Bommasani, R. et al. 2021. On the Opportunities and Risks of Foundation Models. (Aug. 2021).

Insight

# A Unified Life Cycle...

# Foundation Model for Scene Graphs



Foundation Model for Scene Graphs

*Graphs support the relationships between data space resources*

Insight

# Downstream Tasks – Image Captioning + Multimodal Question Answering

# Conclusions

- Common European Data Spaces will require new data management paradigms in order to facilitate large-scale data pooling and sharing

- Semantics have never been more challenging
  - Scale, content space, multimodal, subjectivity, trust, cost, …

- Foundational / Language models and hybrid Neuro-Symbolic approaches provide promising avenues to develop new paradigms
  - Computationally driven
  - Knowledge Graphs + Deep Learning
  - Exact + Approximate best effort results

Insight

# SEMIC Specifications, services and trainings in support of Data Spaces

Alexandra Balahur

# SEMIC Specifications, services and trainings in support of Data Spaces

*Alexandra Balahur*

*DG DIGIT, Interoperability Unit*

*6 October 2023*

# Index

What is the index of this presentation?

European Commission

# Why

Understanding the policy context

# Policy Context

> Europe's digital future will be enabled by a **data-driven economy** and the use of **Artificial Intelligence**, fully respecting EU values and regulations. The public sector also needs to become more data-driven; improve the capability of developing policies and services through the management, sharing and use of data.

> The [European Data Strategy](#) aims to create a single market for data through common European data spaces that benefit from common standards and interoperability protocols.

> [AI legislation](#) and [coordinated plan](#) to foster the development and use of AI in Europe, highlighting the public sector as a trailblazer for using AI.

> The [Interoperable Europe Act](#) complements the EU data and digital policy landscape on data availability and data exchange, from a public sector angle. It implements interoperability by design and fosters the sharing and reuse of interoperable solutions.

# What

Understanding the role and services from DG DIGIT

# DG DIGIT's role and services

DG Informatics is ready to support data spaces with existing assets and services, as well as to establish synergies with stakeholders active in this field to provide a more comprehensive support.

## European Commission DG Informatics

**Interoperability**

**Analytics**

**Trust & Identity**

**Interoperable Europe**

## Support to data spaces

**CATALOGUES & METADATA**

**DATA MODELS**

**TRUSTWORTHY DATA EXCHANGE**

**CONSULTANCY & PILOTING**

DCAT-AP FOR DATA PORTALS IN EUROPE

#GOVERNMENT CORE VOCABULARIES

APPLICATION PROFILES FOR DATA PORTALS IN EUROPE

ebsi

Personal Data Spaces

Linked Data Event Streams

Wikidata Wikibase

Mortgage accepted
eSign contract

eGov Portal
Deliver Documents

Insights Revealed
Data Analytics

**DSSC SERVICES**

## Data spaces

Health    Tourism    Mobility

Manufacturing    Green Deal    Skills

Agriculture    Energy    Public procurement

Finance    Cultural heritage    Languages

## Data Spaces Support Centre

DATA SPACES SUPPORT CENTRE

BDV BIG DATA VALUE ASSOCIATION

FIWARE

INTERNATIONAL DATA SPACES ASSOCIATION

KU LEUVEN

gaia-x

GAILLIMH GALWAY

European Commission

# Focus on Semantic Interoperability

Understanding the role and services from SEMIC supporting data spaces

European Commission

# SEMIC Service offering

SEMIC's goal is to deliver pragmatic support to help build an interoperable Europe.

## Specifications

Publication and maintenance of open and free-to-reuse data models, with regular updates

## Pilots

Developing specific solutions for public administrations to scale up their interoperability maturity

## Toolkit

Provision of an accessible European Toolchain for data extraction, transformation and loading

## Knowledge Hub

Training materials, guidelines and events to foster interoperability and share knowledge of its benefits

European Commission

# SEMIC specifications

SEMIC specifications enable interoperability by:

- Making data transparent and available
- Supporting coherent implementation of laws and policies
- Helping implement cost efficiencies
- Helping digitalisation and harmonising processes

## Core Vocabularies

The cornerstone for semantic interoperability

Core Vocabularies provide a standardised approach for describing key concepts such as locations, businesses, organisations and natural persons.

## Application Profiles

A tailored data model for specific applications

Application Profiles make use of vocabularies for a detailed set of use cases to define mandatory relations, constraints and relationships.

We adopt a balanced approach providing flexibility, customisation and solid indications to ensure a high degree of semantic interoperability.

## Examples

SEMIC Core Vocabularies are reused in various implementations and extensions across domains and countries, from INSPIRE regulation to national databases (e.g. organisation registers) and private initiatives (Smart Data Models from Fiware).

DCAT-AP is a standard model to describe data catalogues. It is used by all open data portals in Europe. We have already made available extensions for statistics, geospatial and base registries.
We are currently supporting the European Health Data Space and the Mobility Data Space to extend DCAT-AP in their respective domains. This will ensure that minimum data descriptions will flow across data spaces.

CPSV-AP is a reusable and extensible data specification used for harmonising the way public services are described in a machine-readable format. Public administrations and service providers use this approach to describe their services and guarantee a level of cross-domain and cross-border interoperability at European, national and local level.



DCAT
↓
DCAT-AP

Extensions
StatDCAT-AP
GeoDCAT-AP
BRegDCAT-AP

CPSV
↓
CPSV-AP
For the provision and description of public services

ADMS
↓
ADMS-AP
For the description of interoperability assets

# Zooming into DCAT-AP

## Objectives of DCAT-AP

Supporting the discovery of/access to (open) data in a cross-border and cross-domain environment, by describing metadata to be harvested across a distributed network of portals. DCAT-AP is based on the Data Catalog Vocabulary (DCAT) from W3C and

- expresses constraints and usages on DCAT properties and classes, and
- includes additional properties and usages of controlled vocabularies

## Domains of applications

**Open data portals** with an extension for statistics and geospatial data.

**Base registries** metadata descriptions

**Machine Learning** with MLDCAT-AP

**Data spaces**
DCAT-AP can make a significant contribution in connecting data catalogues of data spaces. Stakeholders that we have discussed with have shared their interest in it for several European common data spaces. Examples include:
- NAPCORE-Mobility
- HealthDCAT-AP
- …

## SEMIC DCAT-AP Support in 2023

**7 November 2023**

DCAT-AP working group webinar on status and governance

**21 November 2023**

DCAT-AP working group webinar on technical issues

**Ongoing**

Public review DCAT-AP 3.0.0

**December 2023**

Training materials on DCAT-AP, dedicated to Data Spaces

**December 2023**

Guidelines on profiling and extending DCAT-AP

European Commission

# Benefits of the DCAT-AP ecosystem

**Enhances the findability** and accessibility of data

Comes with a decade of **experience** of documenting, maintaining metadata records; sharing through harvesting, etc.

Provides tooling to **validate the implementation** data.

Enables data spaces to make **data catalogues findable**
→ A harvesting network is made possible

Enables data spaces to express their **metadata in a common language**

**Collaborative environment** that allows data spaces to express their needs and additional requirements (specialisations)

European Commission

# SEMIC pilots

SEMIC sets up pilots to showcase the value of new approaches and ecosystems, which can be leveraged across public administrations to scale up their interoperability maturity. Pilots usually involve participants from several Member States and sector-specific DGs co-creating solutions with SEMIC's support.

## Artificial Intelligence

Elevating the application of AI tools

SEMIC wants to become the reference in AI for interoperability and interoperability in AI for public organisations in Europe.

## Personal Data Spaces

Empowering individuals and supporting the data economy

There is a need for more coordination and synergies between Personal Data Spaces implementations to ensure interoperability between existing and potential solutions in this emerging market.

## Linked Data Event Streams

A new data publishing approach

A publishing strategy by which a data provider allows multiple third parties to stay in sync with the latest or historical versions of the data source in a cost-effective manner.

## Wikidata / Wikibase

Supporting community-driven efforts

Tools to enable the co-creation of semantic data models emerging from a community of users.

## Examples

**Linked Data Event Streams**

The European Railways Agency has deployed a solution based on the LDES technology which allowed them to save costs and efficiently transfer only the data that has changed, together with what has changed and when, for historic records.

**Proof-of-concept**

Proof-of-concepts serve as practical examples of how AI can complement and automate existing ways of working, facilitating more efficient work and promoting interoperability within the European Union's stakeholders. These tools demonstrate the potential for AI to improve data processing, analysis, and interoperability. For example, the NLP-based proof-of-concept for automatically tagging web pages.

**Studies and research**

Studies and research conducted under SEMIC are research endeavours aimed at investigating various aspects of AI, interoperability, and advanced technologies. For example, SEMIC is currently creating a conference paper on fine-tuning of Large Language Models in Tourism domain.

European Commission

# Zooming into LDES

## Use cases

**Replicate the data:** An LDES allows all data users to replicate the data from the unique source of truth (i.e. a base registry).

**Stay in-sync:** Stay up-to-date on the latest changes of the event source.

**Grow your collection of objects:** Each time an event happens, a new object will be added to the LDES allowing you to update your data on an event-base.

## Benefits

**No more maintenance hell**
The data owner does not have to maintain all the querying APIs on their side.

**Less chance for replication hell**
There is no need for a local copy of the data source since it easy to connect and stay up-to-date.

**Publishing data in a scalable manner**

Since objects are published each time an event happens, the LDES publication technology is scalable by nature.

## The role of SEMIC in LDES

A Linked Data Event Stream (LDES) is a collection of immutable objects whereby you do not change the data itself but simply add new data record to the stream. For business purposes, it is a publication strategy to share your data. It allows data users to:

- Have up to date data
- Be aware of changes
- Access historic data
- Relate historic data to current data

SEMIC assists Members States and organisations to publish their data in a manner that ensures interoperability, using LDES:

- Aid in the implementation of pilots
- Aid in the development of Strategy and roadmap
- Share knowledge

European Commission

# SEMIC Knowledge Hub

SEMIC acts as an enabler to exchange experiences, good practices and insights. By sharing knowledge we aim to facilitate the development and use of data specifications, as well as to discuss the latest technological trends, and present expert views on semantic interoperability topics.

## Community building

Experience exchange and building consensus

Organisation of webinars, workshops and other events to understand user and market needs, foster experience and good practice exchange, and help reach consensus between stakeholders.

## Studies & other materials

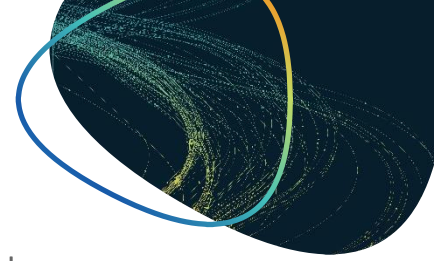Setting trends and providing guidance

Over the years SEMIC has published many reports, studies and other guidance materials revolving around semantic interoperability.

## Training

Providing online learning materials

SEMIC contributes to advancing digital skills in the area of interoperability to support policy, service delivery and impact evaluation.

## Examples

SEMIC has been developing many online training materials that contribute to advancing digital skills. Some of the training developed are particularly relevant to data spaces.

The training developed by SEMIC can be grouped into three main categories:
- Core semantics
- DCAT and DCAT-AP
- Other semantic solutions for public administration

The training materials developed by SEMIC are all publicly available on EU Academy.

European Commission

# Online trainings with relevance to data spaces

The training materials developed by SEMIC are all publicly available on EU Academy.

**CORE SEMANTICS**
Introduction to Core Vocs
Introduction to CPSV-AP
ABR eLearning course


Core Semantics
**Introduction to Core Vocabularies** — Beginner
Public administrations, Semantic experts & IT architects


Core Semantics
**Introduction to the Core Public Service Vocabulary Application Profile (CPSV-AP)** — Intermediate
Analysts who are interested in expanding their CPSV-AP knowledge


Core Semantics
**Access to Base Registries eLearning course** — Intermediate
People who describe data assets in the base registries, public administrations and data stewards

**DCAT AND DCAT-AP**
Introduction
Basic user
Advanced user


DCAT and DCAT-AP
**DCAT and DCAT-AP training: General introduction** — Beginner
People who describe data assets in internal catalogues and/or on EU ODP, and data stewards in EU institutions


DCAT and DCAT-AP
**DCAT and DCAT-AP training: Basic user** — Beginner
Users with basic DCAT knowledge, and participants of the 'DCAT and DCAT-AP training: General introduction'


DCAT and DCAT-AP
**DCAT and DCAT-AP training: Advanced user** — Proficient
Advanced users with working experience in DCAT/DCAT-AP, and users with substantial knowledge

**OTHER SEMANTIC SOLUTIONS FOR PUBLIC ADMINISTRATIONS**
SOLID
Wikibase
LDES


Semantic solutions
**Introduction to SOLID** — Intermediate
People interested in expanding their knowledge of SOLID, SOLID implementers, and computer science students


Semantic solutions
**Wikibase and Semantic MediaWiki for data driven semantics** — Intermediate
Data maintainers and IT professionals, public administrations and policy makers


Semantic solutions
**Publishing data with Linked Data Event Streams: why and how** — Beginner
Managers, data maintainers and developers thoroughly interested on their Linked Data's life cycle

# Where

Reaching out to us through a single entry point

European Commission

# SEMIC Support Centre

Your one-stop-shop for all digital interoperability challenges

## What is it?

- Concrete and direct support to all stakeholders through our helpdesk and GitHub
- Collects good practices and serves as a knowledge hub on interoperability issues

## Who is it for?

- Public officers and technical experts
- Any public organisation in its interoperable journey

# SEMIC Support Centre – Data Spaces



## What can be found on the page?

- Explanation about what data spaces are

- Information about how DG DIGIT supports data spaces

- Specific support provided by DIGIT B

- Contact details to ask for support

# Contact us

Find more information and contact us through a single entry-point:

[https://joinup.ec.europa.eu/collection/semic-support-centre/data-spaces](https://joinup.ec.europa.eu/collection/semic-support-centre/data-spaces)

Thank you !

# Questions & Answers

Please provide your feedback!

Stay up-to-date on our
**2023 activities!**

data.
europa
academy

# Join our next webinars!



Workshop

**How to use open data for your research?**

data. europa academy

19 October 2023
10.00 – 12.00 CET



WEBINAR

**Open Data Maturity 2022: Diving deeper into the impact dimension**

data. europa academy

27 October 2023

10.00 — 11.30 CET

# Fill in our user survey

**Sign up for the newsletter:** data.europa.eu/newsletter
**Follow us on social media:**

EU_opendata

Publications Office of the European Union

data.europa.eu

**Thank you**

data.
europa
academy